# Large scale long-tailed product recognition system at Alibaba

Xiangzeng Zhou, Pan Pan, Yun Zheng, Yinghui Xu, Rong Jin
Machine Intelligence Technology Lab, Damo Academy
Alibaba Group, Hangzhou, China
xiangzeng.zxz,panpan.pp,zhengyun.zy@alibaba-inc.com
renji.xyh@taobao.com,jinrong.jr@alibaba-inc.com

## ABSTRACT

A practical large scale product recognition system suffers from the phenomenon of long-tailed imbalanced training data under the E-commercial circumstance at Alibaba. Besides product images at Alibaba, plenty of image related side information (e.g. title, tags) reveal rich semantic information about images. Prior works mainly focus on addressing the long tail problem in visual perspective only, but lack of consideration of leveraging the side information. In this paper, we present a novel side information based large scale visual recognition co-training (SICoT) system to deal with the long tail problem by leveraging the image related side information. In the proposed co-training system, we firstly introduce a bilinear word attention module aiming to construct a semantic embedding over the noisy side information. A visual feature and semantic embedding co-training scheme is then designed to transfer knowledge from classes with abundant training data (head classes) to classes with few training data (tail classes) in an end-to-end fashion. Extensive experiments on four challenging large scale datasets, whose numbers of classes range from one thousand to one million, demonstrate the scalable effectiveness of the proposed SICoT system in alleviating the long tail problem. In the visual search platform Pailitao[1] at Alibaba, we settle a practical large scale product recognition application driven by the proposed SICoT system, and achieve a significant gain of unique visitor (UV) conversion rate.

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

## KEYWORDS

product recognition, long-tailed, attention, side information, co-training

**ACM Reference Format:**
Xiangzeng Zhou, Pan Pan, Yun Zheng, Yinghui Xu, Rong Jin. 2020. Large scale long-tailed product recognition system at Alibaba. In *Proceedings of the 29th ACM International Conference on Information and Knowledge*
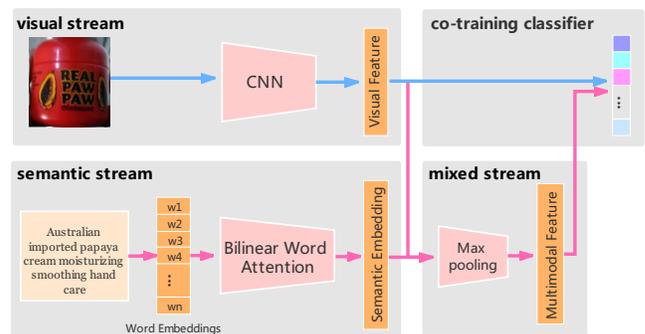
---

[1]http://www.pailitao.com

**Figure 1: The overall architecture of our proposed side information based co-training (SICoT) system. The system contains four streams: (i) The visual stream is a visual feature extractor using a convolutional neural network. (ii) The semantic stream, which consists of a word2vec module and a proposed bilinear word attention module, aims to learn a semantic embedding from the noisy side information. (iii) The mixed stream takes charge of generating a multimodal feature. (iv) The shared classifier is co-trained by the visual stream and the mixed stream.**

## 1 INTRODUCTION

Recent years have witnessed the remarkable progresses of wielding deep learning models in visual recognition task. With the aid of deep learning techniques, nowadays it is practicable to establish an industrial large scale visual recognition application based on huge volume of image data. Compared to the quantity of products and image data under the E-commercial circumstance at Alibaba, many popular so-called large scale visual recognition datasets, like ImageNet [6], WebVision2.0 [16], iMaterialist Product 2019 [13] and Open Images V4 [15], appear to be relatively small scale. Even though some datasets have reached several millions of training image data, the quantity of categories only ranges from hundreds to thousands.

On the basis of abundant image data of great value in the E-commercial scenario and powerful computing resources at Alibaba, it is still great challenging to establish a truly practical large scale visual recognition application. And these challenges are roughly reflected in following three aspects:

**Enormous quantity of classes and images:** There are about tens of million of daily active products and billions of image data

in the marketplace of Alibaba, which covers categories of clothing, shoe, bag, cosmetic, drink, snack and toy in general. The huge quantity of product classes and images brings difficulties to both the training and deployment of a large scale visual recognition model. For example, when the number of classes reaches about one million, the size of the last fully connected (FC) layer will exceed the maximum memory of a single block of Nvidia-V100-32G GPU. This requires a new training paradigm capable of training a huge FC layer in a distributed manner.

**Scarce and noisy annotation:** Unlike those well compiled small scale datasets (e.g. ImageNet [6]), it is impracticable to manually annotate the daily growing enormous quantity of image data in the E-commercial scenario. Training images with high-quality annotations are scarce and insufficient to build a practical large scale visual recognition system. Although there are lots of annotations provided by sellers or customers in the marketplace of Alibaba, it is of great difficulty to use these weakly and noisy annotations to assist in training a visual recognition model.

**Long-tailed distribution of training data:** The phenomenon of long-tailed imbalanced training data naturally occurs under the E-commercial circumstance. In the marketplace of Alibaba, a large amount of new arrival products emerge everyday, meanwhile, quite a large portion of products are low sales or even zero sale. The difficulty to acquire sufficient training images for these products challenges the performance of a large scale visual recognition application.

It is known that, without any special treatment of the classes with insufficient training data (tail classes), the classification boundary of a recognition model inclines toward those classes with abundant training data (head classes). At Alibaba, there are abundant image related data or side information, such as short titles and long text descriptions, coming from various sellers or customers. These side information that containing rich, yet weak and noisy annotations reveal underlying similarity among images and classes from a different perspective.

However, prior works [2, 7–9, 14, 20, 24, 31, 33] mainly focus on addressing the long tail problem in visual perspective only, but lack of consideration of leveraging the side information. On the other hand, most works [4, 12, 19] take advantage of the side information as a kind of weakly supervision in general, and are not meant to address the long tail problem. Inspired by the work of transferring knowledge or borrowing training examples between similar classes [17] in detection task, we attempt to address the problem of long-tailed distributed training data in the task of large scale product recognition by leveraging the noisy side information in this paper. Considering the following two facts observed on the data of the marketplace at Alibaba, the usage of image related side information has great potential to alleviate the long tail problem. a) Unlike the extreme imbalanced distribution of image data, the distribution of words from image related titles is relatively much balanced. b) About 12% words out of the entire vocabulary are shared between the head classes and the tail classes.

In this paper, we propose a novel side information based visual recognition co-training (SICoT) system, as shown in Fig. 1, which aims to deal with the long tail problem in a large scale product classification task. Moreover, we have launched a SKU level product recognition service driven by the proposed SICoT system in the
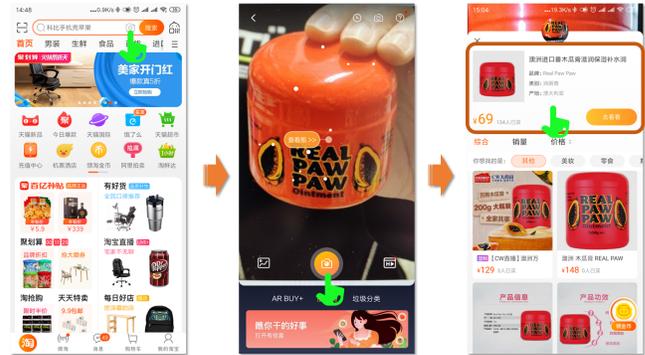


**Figure 2: The scenario of large scale SKU level product recognition service on Pailitao at Alibaba: by taking a picture, Pailitao identifies the product with its short title and several associated tags in real time shown at the top of the page.**

visual search platform Pailitao [1, 32] at Alibaba, as shown in Fig. 2. We conclude our contributions as following:

1) We introduce a bilinear word attention module to distinguish important words from the noisy side information of image short titles, followed by constructing a semantic embedding as a kind of distilled knowledge of the side information.

2) Considering the long tail problem in a large scale product recognition task, we propose a novel visual feature and semantic embedding co-training (SICoT) system to help transfer knowledge from the head classes to the tail classes in an end-to-end way. This co-training system aims to perform transfer learning across the head and the tail classes by deeply involving the side information in both feature learning and classifier training.

3) Extensive experiments on both open large scale datasets and our organized huge scale SKU level product datasets demonstrate the scalable effectiveness of the proposed side information based co-training system in relieving the long tail problem.

## 2 RELATED WORK

Taking the issue of long-tailed distributed training data into account, it still remains very challenging problems on how to establish a practical large scale product recognition system in the E-commercial scenario at Alibaba. Some prior works about addressing the long tail problem, taking advantaging of side information and large scale product recognition are roughly summarised in the following aspects.

**Imbalanced learning:** To alleviate the problem of long-tailed distributed training data, many traditional approaches have been extensively studied in the past [2, 7–9, 14, 20, 24, 31, 33]. Re-sampling methods [9, 20] aim to balance the numbers of training samples between multiple classes by under-sampling the head classes or over-sampling the tail classes. However, these methods often lead to removal of important samples or introduction of meaningless duplicated samples. Cost-sensitive methods [24, 33] try to make the standard classifiers more sensitive to the head classes by imposing higher misclassification cost to the head classes than to the tail classes. Most recently, deep neural networks are widely applied to performing imbalanced learning [2, 7, 8, 14, 31, 33]. Apart from the
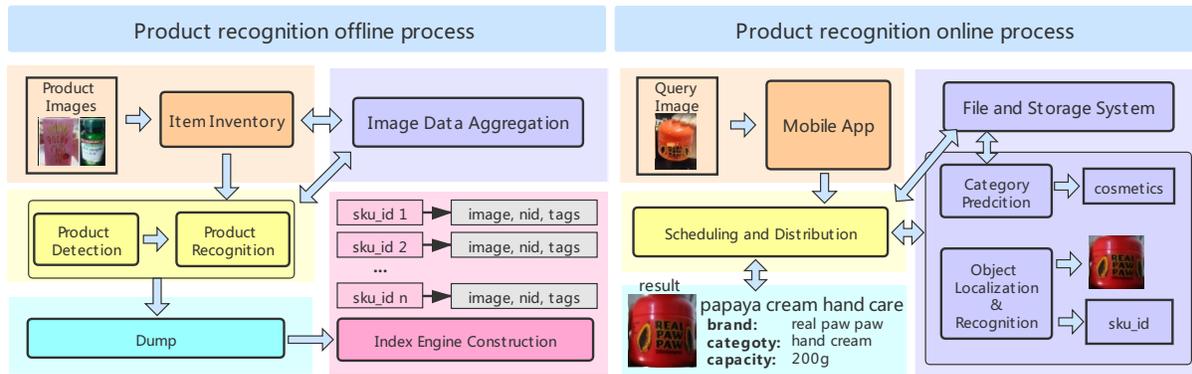
**Figure 3: Overview of the overall product recognition architecture settled in Pailitao.**

conventional cross entropy loss, several new objective loss function, like range loss [31], class rectification loss [7] and cluster-based large margin local embedding [8] are proposed to address the long tail problem by rectifying the classification boundaries dominated by the head classes. The works mentioned above devote major effort to addressing the long tail problem merely from the visual aspect. Considering that the side information can provide rich semantic information about image, in this paper the side information are involved in helping tackle the long tail problem.

**Weakly supervised learning:** In many popular datasets and practical scenarios, lots of auxiliary data or side information associated with images is provided, such as image titles and long text descriptions of in WebVision2.0 [16], wordnet in ImageNet [6] and so on. These side information normally come from heterogeneous data sources via web search, and naturally contain a lot of noise. In most of prior works, the noisy side information is mainly taken as a kind of weakly supervision in coordination with other kinds of learning tasks [4, 12, 19]. However, especially in a classification task, taking advantage of the side information as supervision in the one-hot fashion may not fully exploit the knowledge and be somewhat shallow. Besides, the usage of the side information is these works are not meant to address the long tail problem. In our proposed SICoT system, we explore a novel fashion to handle the long tail problem by taking advantage of the side information.

**Knowledge distill and transfer:** The basic principle of transfer learning are also introduced to transfer knowledge or even borrow training data from the head classes to the tail classes [17, 22, 34]. A series of works [27, 28, 30] present a learning using privileged information (LUPI) framework to transfer knowledge (e.g. similarity or margin) across multiple models which usually learned in different modalities. A teacher-teach-student scheme presented in the LUPI framework provides a relatively deeper way to use the side information to assist the original visual recognition task. Extending to this, [18, 29] unify the LUPI framework and the knowledge distillation paradigm [11] into a generalized distillation framework. These works provide a possible way to use the side information as a teacher model to help a visual recognition task (student model) in a teacher-teach-student or distillation scheme. However, the teacher-teach-student scheme requires the teacher model to be learned beforehand. And this two-stage approach is hard to be optimized

globally. This scheme also implicitly demands the teacher model to be more powerful than the student model. In our proposed SICoT system, the side information participate into the visual classification task in conjunction with the visual feature in an end-to-end way. Moreover, no prior assumptions and constraints are made on the side information.

**Large scale product recognition:** Taking the applied techniques into account, Trax [3, 26] and MalongTech [25] have devoted their efforts to establish practical large scale visual recognition applications, especially the SKU (stock keeping unit [2]) level product visual recognition. MalongTech hosts a SKU level product dataset iMaterialist product 2019 and a corresponding competition in conjunction with FGVC6 workshop of CVPR2019 [13]. However, the iMaterialist product 2019 dataset covers only two thousands of product SKUs, and only provides image data without any side information for extending research. At Alibaba, in order to establish a scalable product recognition system, we organize a large scale SKU level product dataset that consists of about 60 million real-shot product images covering 1 million SKUs with the aid of the visual search engine Pailitao [1, 32]. Apart from images, the dataset contains abundant yet noisy side information (e.g. image titles, long descriptions and tags) provided by various sellers or customers in the marketplace of Alibaba. On the basis of the large scale SKU level product dataset, we settle a 30 million products recognition service driven by the proposed SICoT system in Pailitao. To our knowledge, this is the largest scale product recognition application in the E-commercial scenario so far.

## 3 APPROACH

In this section, we elaborate our proposed side information based co-training (SICoT) system in a large scale visual recognition task over long-tailed distributed training data. In Sec. 3.1, we firstly illustrate the overall product recognition process on the basis of the visual search service Pailitao [1, 32]. Considering that the side information of image titles from heterogeneous resources are often noisy, we then propose a bilinear word attention network in Sec. 3.2 to distinguish the important words from the noisy side information. Subsequently, a detailed illustration of the side information based

---

[2]https://en.wikipedia.org/wiki/Stock_keeping_unit

show dog english champ
setter pedigree gundog

autumn mountain
englishsetter

Coca Cola Soft Drink Soda Drink
330ml*24 cans Classical Old/New
Package Random Shipping

Coca Cola Sprite Refreshing
Lemon Flavor Soft Drink 330ml*8
Cans Slim Can Free Shipping

**Figure 4: Two images of English setters from the Webvision2.0 [16] and two images of products from the marketplace of Alibaba. Along with each image, a short text description is also provided. These text descriptions are naturally noisy.**

co-training system (SICoT) is given in Sec. 3.3. The SICoT system aims to leverage visual features and semantic embeddings to help transfer knowledge from the head classes to the tail classes in an end-to-end fashion.

## 3.1 Product Recognition Architecture

The entire product recognition architecture, as shown in Fig. 3, comprises an offline process and an online process, that following the present visual search architecture of Pailitao [1, 32] in general. The offline process mainly refers to the daily process of building product index using the proposed SICoT product recognition system. Unlike the index engine in the visual search service, the product index stores the SKU ids of products and corresponding product images, titles and tags for online retrieval. In the online process, the core function is a real time product recognition service in charge of predicting a SKU id for each query image. For the other modules in this architecture, we simply reuse the design of the visual search service, like category prediction and object localization. Once a query image is successfully recognized by the online service, a predicted SKU id will be obtained. By retrieving the index engine using the SKU id, the corresponding product image, title and tags will be obtained and presented to the customer, as shown in Fig. 2.

## 3.2 Bilinear word attention network

In many practical scenarios, besides images lots of related side information (e.g. image titles) can be obtained. These image related side information may reveal underlying similarity among images and classes, so as to it has great potential to improve a visual recognition task. In order to process both the visual information of images and the image related side information in an unified framework, we propose a bilinear word attention network, as shown in Fig. 5, to learn a semantic embedding from the noisy side information.

A conventional tokenization is firstly conducted on the side information of image titles, followed by using a word2vec model to generate a word embedding for each word after tokenization. In most of natural language processing (NLP) related tasks, such as language translation and image captioning, the order of words in a title matters in general. However, in our experimental observations, the order of words is less important and even harmful, especially when the side information is highly noisy. Here, we simply use an average pooling operation to generate a global embedding from the word embeddings instead of using a sequential model (e.g. recurrent neural networks). As for each word embedding, this global embedding can be regarded as a global context without consideration of the order of words in the side information.

Noise and meaningless words naturally occur in the image related side information due to the heterogeneous resources in both the marketplace of Alibaba and many open datasets. For example, as shown in Fig. 4, both the text descriptions of two dogs from Webvision2.0 and the titles of two soft drinks from the marketplace contain several words less relevant to the image content (e.g. *show*, *autumn*, *free shipping*). Considering this issue, we propose a soft attention sub-network to evaluate the importance of each word in the side information. A bilinear operation between the global context and all word embeddings is introduced to generate a second order feature map. An average pooling operation and a nonlinear transformation are then carried out above the feature map to output a word attention vector. A semantic embedding of the entire side information is achieved by the weighted sum of the word attention vector and all word embeddings at final.

## 3.3 Side information based co-training system

Given the bilinear word attention based semantic embedding presented in Sec. 3.2, we propose a visual feature and semantic embedding co-training scheme in this section. Due to the absence of image related side information once a recognition model has been deployed, the proposed co-training scheme is only involved in the training stage. The scheme is designed to take the semantic embedding from the side information as a bridge to transfer knowledge from the head classes to the tail classes. Unlike the teacher-teach-student paradigm used in LUPI [27], our approach makes no prior assumption about models, i.e. a teacher model should be more powerful than a student model. In fact, in our experimental observations, it is often inadequate to carry out a satisfactory classification by using the image related side information only, especially when the number of class are huge.

As illustrated in Fig. 1, we show the overall architecture of our proposed co-training scheme that consisting of three streams, i.e. a visual stream, a semantic stream and a mixed stream. The visual stream is a conventional visual recognition pipeline, which comprises a common convolutional neural network as a feature extractor and a plain softmax classifier optimized by a cross entropy loss. The visual stream is the target task that we attempt to improve. The semantic stream is simply the bilinear word attention based semantic embedding sub-network, in which the word embeddings are required initialized from a pretrained model like Word2vec [21]. Noted that it is assured that the visual feature and the semantic embedding take a same dimension through a deliberate design of
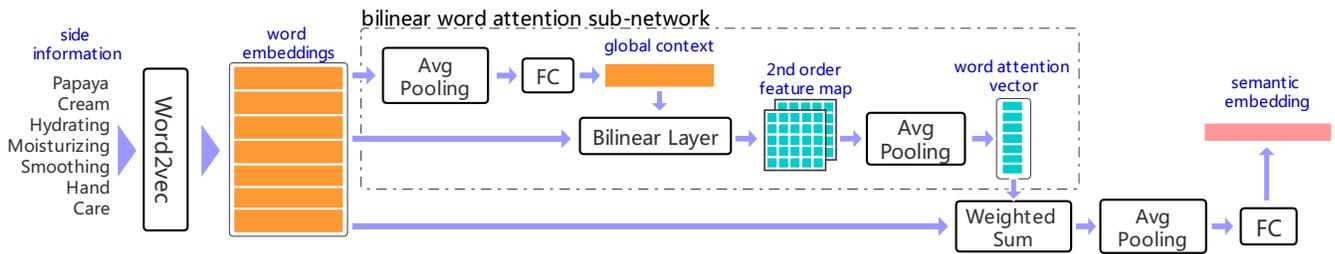
**Figure 5: The semantic embedding network mainly consisting of a word2vec module and the proposed bilinear word attention module.**

the network. A mixed or multi-modal feature is then achieved by a max-pooling operation over the visual features and the semantic embeddings. Unlike a conventional multi-task framework that consisting of multiple learning tasks driven by different objectives, the proposed co-training scheme makes both the visual features and the semantic embeddings to be learnt driven by a same classification task. As shown in Fig. 1, the visual feature $x^v$ and the multi-modal feature $x^m$ are followed by a shared co-training classifier optimized with the objective as Equ. 1, in which the $\hat{y}_i^v$ and $\hat{y}_i^m$ are the classifier output of the visual feature $x_i^v$ and the semantic embedding $x_i^m$, respectively.

$$Loss = -\frac{1}{N}(\sum_{i=1}^{N} y_i \cdot \log(\hat{y}_i^v)) - \lambda \cdot \frac{1}{N}(\sum_{i=1}^{N} y_i \cdot \log(\hat{y}_i^m)) \quad (1)$$

In the proposed training scheme, the shared classifier are learnt in a co-training fashion, in which the classification boundaries are directly affected both visually and semantically. It is well known that the classification boundaries in a classification task with long-tailed imbalanced training data are easily dominated by the head classes. The proposed co-training scheme may rectify the skewed classification boundaries by introducing the semantic knowledge into the classification training. Furthermore, it can be observed in Fig. 1 that the visual and semantic streams are tangled in not only the classifier training part, but also the feature learning part via the gradient backward procedure. Compared to the methods of taking the side information as weakly supervision and the two stage teacher-teach-student paradigm in the LUPI framework, our proposed co-training scheme provides a deeper way to tangle the visual and semantic knowledge in an end-to-end manner.

## 4 EXPERIMENTS

In this section, we evaluate our proposed side information based co-training approach on four large scale datasets that exhibiting long-tailed distribution, i.e. iMaterialist Product 2019 [13], Webvision2.0 selected 1k [16] and our proposed 34k and 1M SKU level product datasets, and demonstrate the positive effect of our approach on the long tail problem. We also report a relative gain of unique visitor (UV) conversion rate after settling our approach to the visual search application Pailitao [1] at Alibaba.

### 4.1 Datasets preparation

The WebVision2.0 [16] dataset is designed to facilitate the research on learning visual representation from noisy web data and it is attached with Google and Flickr retrieval results that containing titles and detailed text descriptions. In this paper, we extract the title of each image as the side information for experiments. We preprocess these side information by discarding the punctuation marks and 5 percent of very frequent and very infrequent words, and throwing away those images without any side information. Eventually, one thousand classes are randomly picked out from the entire 5 thousand of classes in Webvision2.0 [16]. iMaterialist Product 2019 [13] hosted by MalongTech [25] is a SKU level product dataset that consisting of 2019 product SKUs. Because there is no any image related side information provided by Materialist Product 2019, we take each image as an input query of the image search engine Pailitao [1] and collect the title of the top1 search result as the side information. Apart from the two open datasets, we also provide two SKU level product datasets in this paper, a facial cream and mask product dataset that containing 34 thousands classes and a huge scale product dataset that containing one million classes. The two proposed datasets come from the marketplace of Alibaba. An overview of statistics about the four datasets are illustrated in Table 1. And the long-tailed distribution of training data are shown in Fig. 6, in which the boundaries between the head and tail classes on training data are experimentally determined as 200, 5000, 100 and 200, respectively. About the testing sets in the four datasets, each class contains approximately equal number of images for a fair evaluation.

### 4.2 Evaluation metrics

Considering that many classes are conceptually overlapped and ambiguous, especially in the Webvision2.0 and the SKU level product datasets, we report the results of $top1$ and $top3$ predicted labels for the evaluation of recognition performance. In addition, the evaluation of overall $top1$ and $top3$ accuracies are also conducted over the head and tail classes, respectively, to demonstrate the effect on the long tail problem.

### 4.3 Implementation details

In the proposed co-training scheme, we use a Resnet-50 [10] initialized from a ImageNet pretrained model as the backbone convolutional neural network (CNN) in the visual stream. As shown
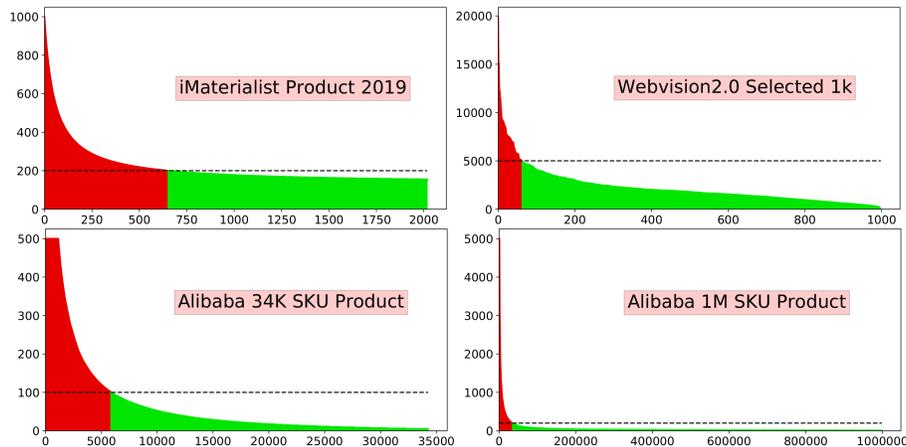
**Figure 6: Long-tailed distribution of training set of the four datasets. The red and green parts in each panel represent the head and tail classes, respectively.**

|  | Class | Trainset | Testset | Vocab |
|---|---|---|---|---|
| iMaterialist Product 2019 | 2019 | 440438 | 9986 | 175208 |
| Webvision2.0 Selected 1k | 1000 | 2230968 | 59040 | 162541 |
| Alibaba 34K SKU Product | 34258 | 2305853 | 171290 | 19762 |
| Alibaba 1M SKU Product | 998131 | 55776960 | 7675690 | 345656 |

**Table 1: Statistics of the four datasets, in which the "Vocab" means the number of words extracted from the side information of image titles after a tokenization process.**

in Fig. 1, the union of the visual stream and co-training classifier in the architecture represents the baseline that carrying out a conventional classifier using image data only. The word embeddings in the semantic stream are initialized using the *word_embedding*() API of Alibaba NLP toolbox (AliNLP) [5]. The word embeddings of AliNLP are trained using the product titles from the marketplace of Alibaba, and well performs on many natural language processing tasks under the E-commercial circumstance at Alibaba.

When training on such a large dataset as our proposed product dataset that containing one million classes, the size of the last fully connected (FC) layer will be larger than the memory size of a single block of GPU (e.g. Nvidia V100 32G). The proposed co-training system is implemented in a hybrid parallel training framework [23], in which the last FC layer is divided and sent to $M$ GPUs for distributed training. The training is carried out in a distributed computing platform of Alibaba with 60 blocks of Nvidia P100 GPUs. For a fair comparison, the baseline and the co-training approach are trained using a stochastic gradient descent (SGD) optimizer with a same learning configuration of a batch size 256, an initial learning rate 0.1 and a step decay policy of step 1, gamma 0.8.

## 4.4 Experimental results

*4.4.1 Overall classification accuracy.* As illustrated in Table 2, we report the $top1$ and $top3$ (enclosed in parentheses) classification

accuracies on the four datasets. It is clearly observed that our proposed co-training approach using the side information of image titles shows performance improvements against the baselines with significant $top1$ (and $top3$) accuracies gain by 1.01%(1.68%) in iMaterialist Product 2019, 3.27%(1.91%) in Webvision2.0 Selected 1k, 1.87%(1.85%) in Alibaba 34K SKU Product and 2.00% (0.86%) in Alibaba 1M SKU Product. When the number of classes ranging from 1 thousands to 1 millions, our approach achieves consistent and significant accuracy gains all the time. The consistent improvement of classification performance demonstrates an attractive scalability for setting out practical visual recognition applications.

*4.4.2 Effect of the bilinear word attention based semantic embedding.* For a qualitative evaluation of the proposed bilinear word attention based semantic embedding, as illustrated in Fig. 7, we visualize the learned word attention vectors of several title samples. The titles of the four examples are demonstrated in original Chinese and translated English at the same time. These words are displayed in the descending order of the corresponding value in the attention vector. In each example, three most important words are highlighted in green and three most unimportant words are highlighted in red. It is observed that the words with larger attention weights generally refer to product names and brand names, and the words with smaller attention weights are often less helpful to distinguish the product from other products. For example, the promotion words of "discount" and "new arrival" are irrelevant to identifying a product.

|  | Baseline | Co-training | Gain |
|---|---|---|---|
| iMaterialist Product 2019 | 57.29 (83.77) | 58.30 (85.45) | **1.01 (1.68)** |
| Webvision2.0 Selected 1k | 62.68 (79.02) | 65.95 (80.93) | **3.27 (1.91)** |
| Alibaba 34K SKU Product | 51.60 (77.04) | 53.47 (78.89) | **1.87 (1.85)** |
| Alibaba 1M SKU Product | 88.25 (97.13) | 90.25 (97.99) | **2.00 (0.86)** |

**Table 2: The $top1$ and $top3$ (in parentheses) classification accuracies and performance improvements compared to the baselines on the four datasets.**



**Figure 7: Visualization of learned attention word vectors on four examples. Each example consists of a image, an original image title in Chinese and a series of words with corresponding attention value in both Chinese and English. These words are displayed in the descending order of the corresponding value in the attention vector. In each example, three most important words are highlighted in green and three most unimportant words are highlighted in red.**

Meanwhile, as shown in Tab. 3, we compare our proposed bilinear word attention based embedding with several other methods on the proposed Alibaba 1M SKU product dataset. The method of Alinlp embeddings mean pooling simply takes the mean of the Alinlp word embeddings as the final title embedding. The bi-LSTM method takes advantage of a bidirectional long short-term memory (LSTM) module to generate a final title embedding. The bi-LSTM attention method introduces a self attention mechanism into the bi-LSTM method. It is observed that our proposed bilinear word attention based embedding outperforms the other four methods. The two bi-LSTM based methods take the order of words into consideration by using a sequential model bi-LSTM to describe the context. However, the inferior performance of the two methods indicate that the order of words is less helpful to carrying out a classification task. Compared to the Alinlp embeddings mean pooling method, our proposed attention method has a positive effect on boosting the performance of classification by using the bilinear word attention model.

| Methods | Top1 Accuracy | Gain |
|---|---|---|
| baseline | 88.25 | - |
| Alinlp embeddings mean pooling | 89.14 | **0.89** |
| bi-LSTM | 89.13 | **0.88** |
| bi-LSTM attention | 88.62 | **0.37** |
| bilinear word attention | 89.41 | **1.16** |

**Table 3: Comparison of our proposed bilinear word attention based embedding with other four methods on the proposed Alibaba 1M SKU Product dataset.**

*4.4.3 Effect on long-tailed distribution of training data.* As shown in Table 4, we illustrate the effect of our proposed co-training scheme on the problem of long-tailed distributed training data. Compared with the baseline on the four datasets, we report the averaged top1 and top3 classification accuracies of the head and tail classes, respectively. It is observed that our approach achieves improvement

|  |  | #Training Samples | Baseline | Co-training | Gain |
|---|---|---|---|---|---|
| iMaterialist Product 2019 | Head | 200 - max | 57.38 (83.49) | 58.37 (85.08) | **0.99 (1.59)** |
|  | Tail | 1 - 200 | 37.21 (81.40) | 41.86 (83.72) | **4.56 (2.32)** |
| Webvision2.0 Selected 1k | Head | 5000 - max | 55.18 (72.91) | 55.86 (72.55) | **0.68** (**-0.36**) |
|  | Tail | 1 - 5000 | 63.04 (79.26) | 66.23 (81.11) | **3.19 (1.85)** |
| Alibaba 34K SKU Product | Head | 100 - max | 62.98 (84.94) | 63.54 (85.51) | **0.56 (0.57)** |
|  | Tail | 1 - 100 | 49.24 (75.40) | 51.38 (77.51) | **2.14 (2.11)** |
| Alibaba 1M SKU Product | Head | 200 - max | 95.09 (99.53) | 95.61 (99.62) | **0.52 (0.09)** |
|  | Tail | 1 - 200 | 86.87 (96.64) | 89.16 (97.66) | **2.29 (1.02)** |

**Table 4: The averaged** $top1$ **and** $top3$ **(in parentheses) classification accuracies and performance improvements on the head classes and the tail classes, repetitively.**

of the top1 and top3 classification accuracies in both the head classes and the tail classes. Noted that the performance gains in the head classes are more significant than the gains in the tail classes. This phenomenon indicates that our approach improves performance of the tail classes while does not harm the performance of the head classes. In fact, our approach often can slightly improve the performance of the head classes at the same time. This mainly owes to that the procedure of knowledge transfer is bidirectional in the proposed co-training system, and the knowledge extracted from the tail classes is also beneficial to the training of the head classes.

*4.4.4 Application in Pailitao.* Pailitao [1] is a visual search application which aims to assisting customers to find the same or similar products by a mobile phone camera shot image. Pailitao is still experiencing swift growth of daily active users (DAU). In Pailitao, the unique visitor (UV) conversion rate is a common measurement which is calculated as Equ. 2.

$$\text{UV conversion rate} = \frac{\text{number of trading UV}}{\text{number of visiting UV}} \quad (2)$$

To further improve the UV conversion rate of Pailitao, a huge scale SKU level product recognition service is settled in Pailitao as an upgrade. As shown in Fig. 2, the recognition result of a query is displayed at the top panel of the page in conjunction with the results of visual search. The panel includes a clickable image, a short title and several tags of the recognized product. The product recognition service is constructed by using the proposed side information based co-training system. Compared to the original version of Pailitao, there is a relative 3.1 percent gain of daily UV conversion rate after this upgrade.

## 5 CONCLUSION

The phenomenon of long-tailed imbalanced training data naturally occurs under the E-commercial circumstance and challenges the performance of a large scale visual recognition task. In this paper we address the problem of long-tailed distributed training data by exploring a side information based visual recognition co-training (SICoT) system. We firstly introduce a bilinear word attention sub-network to distinguish important words from the noisy side information, followed by generating a semantic embedding as the distilled knowledge of the side information. An end-to-end visual feature and semantic embedding co-training system is then proposed to help transfer knowledge from the head classes to the tail

classes. Experimental results on four large scale datasets demonstrate the effectiveness of the proposed approach. Our approach improves the performance of the tail classes without any harm to the head classes. Moreover, a SICoT driven product visual recognition service is settled in Pailitao and achieves a significant gain of unique visitor conversion rate. Our approach has shown good scalability ranging from medium to large scale datasets, and it is of great value for establishing industrial visual recognition applications.

## REFERENCES

[1] Alibaba. [n.d.]. PaiLiTao. http://www.pailitao.com.
[2] Cristiano L Castro and Antônio P Braga. 2013. Novel cost-sensitive approach to improve the multilayer perceptron performance on imbalanced data. *IEEE transactions on neural networks and learning systems* 24, 6 (2013), 888–899.
[3] Daniel Shimon Cohen, Yair Adato, and Dolev Pomeranz. 2019. Method and a system for object recognition. US Patent 10,402,777.
[4] Charles Corbiere, Hedi Ben-Younes, Alexandre Rame, and Charles Ollion. 2017. Leveraging weakly annotated data for fashion image retrieval and label prediction. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2268–2274.
[5] Alibaba DAMO. [n.d.]. AliNLP. https://damo.alibaba.com/labs/language-technology?lang=en.
[6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.
[7] Qi Dong, Shaogang Gong, and Xiatian Zhu. 2017. Class rectification hard mining for imbalanced deep learning. In *Proceedings of the IEEE International Conference on Computer Vision*. 1851–1860.
[8] Q. Dong, S. Gong, and X. Zhu. 2019. Imbalanced Deep Learning by Minority Class Incremental Rectification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 6 (June 2019), 1367–1381.
[9] E. A. Garcia and H. He. 2009. Learning from Imbalanced Data. *IEEE Transactions on Knowledge and Data Engineering* 21, 09 (sep 2009), 1263–1284. https://doi.org/10.1109/TKDE.2008.239
[10] 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.
[11] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).
[12] Armand Joulin, Laurens van der Maaten, Allan Jabri, and Nicolas Vasilache. 2016. Learning visual features from large weakly supervised data. In *European Conference on Computer Vision*. Springer, 67–84.
[13] kaggle. [n.d.]. iMaterialist Challenge on Product Recognition. https://www.kaggle.com/c/imaterialist-product-2019.
[14] Salman H Khan, Munawar Hayat, Mohammed Bennamoun, Ferdous A Sohel, and Roberto Togneri. 2017. Cost-sensitive learning of deep feature representations from imbalanced data. *IEEE transactions on neural networks and learning systems* 29, 8 (2017), 3573–3587.
[15] 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).
[16] Wen Li, Limin Wang, Wei Li, Eirikur Agustsson, Jesse Berent, Abhinav Gupta, Rahul Sukthankar, and Luc Van Gool. 2017. WebVision Challenge: Visual Learning and Understanding With Web Data. *Arxiv Preprint* (2017).

[17] Joseph J Lim, Russ R Salakhutdinov, and Antonio Torralba. 2011. Transfer learning by borrowing examples for multiclass object detection. In *Advances in neural information processing systems*. 118–126.

[18] David Lopez-Paz, Léon Bottou, Bernhard Schölkopf, and Vladimir Vapnik. 2015. Unifying distillation and privileged information. arXiv:1511.03643 [stat.ML]

[19] Yunshan Ma, Xun Yang, Lizi Liao, Yixin Cao, and Tat-Seng Chua. 2019. Who, Where, and What to Wear? Extracting Fashion Knowledge from Social Media. In *Proceedings of the 27th ACM International Conference on Multimedia*. 257–265.

[20] T. Maciejewski and J. Stefanowski. 2011. Local neighbourhood extension of SMOTE for mining imbalanced data. In *2011 IEEE Symposium on Computational Intelligence and Data Mining (CIDM)*. 104–111. https://doi.org/10.1109/CIDM.2011.5949434

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

[22] Ruslan Salakhutdinov, Antonio Torralba, and Josh Tenenbaum. 2011. Learning to share visual appearance for multiclass object detection. In *CVPR 2011*. IEEE, 1481–1488.

[23] Liuyihan Song, Pan Pan, Kang Zhao, Hao Yang, Yiming Chen, Yingya Zhang, Yinghui Xu, and Rong Jin. 2020. Large-Scale Training System for 100-Million Classification at Alibaba (under review). In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery &#38; Data Mining* (San Diego, California USA) *(KDD '20)*. ACM, New York, NY, USA, 993–1001.

[24] Y. Tang, Y. Zhang, N. V. Chawla, and S. Krasser. 2009. SVMs Modeling for Highly Imbalanced Classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 39, 1 (Feb 2009), 281–288. https://doi.org/10.1109/TSMCB.2008.2002909

[25] Malong Technologies. [n.d.]. ProductAI. https://www.productai.com/home.

[26] Trax. [n.d.]. Traxretail. https://traxretail.com/.

[27] Vladimir Vapnik and Rauf Izmailov. 2015. Learning Using Privileged Information: Similarity Control and Knowledge Transfer. *Journal of Machine Learning Research* 16, 61 (2015), 2023–2049. http://jmlr.org/papers/v16/vapnik15b.html

[28] Vladimir Vapnik and Akshay Vashist. 2009. A new learning paradigm: Learning using privileged information. *Neural Networks* 22, 5 (2009), 544 – 557. https://doi.org/10.1016/j.neunet.2009.06.042 Advances in Neural Networks Research: IJCNN2009.

[29] Weiran Wang. 2019. Everything old is new again: A multi-view learning approach to learning using privileged information and distillation. *arXiv preprint arXiv:1903.03694* (2019).

[30] X. Yang, M. Wang, and D. Tao. 2018. Person Re-Identification With Metric Learning Using Privileged Information. *IEEE Transactions on Image Processing* 27, 2 (Feb 2018), 791–805.

[31] X. Zhang, Z. Fang, Y. Wen, Z. Li, and Y. Qiao. 2017. Range Loss for Deep Face Recognition with Long-Tailed Training Data. In *2017 IEEE International Conference on Computer Vision (ICCV)*. 5419–5428. https://doi.org/10.1109/ICCV.2017.578

[32] Yanhao Zhang, Pan Pan, Yun Zheng, Kang Zhao, Yingya Zhang, Xiaofeng Ren, and Rong Jin. 2018. Visual Search at Alibaba. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery &#38; Data Mining* (London, United Kingdom) *(KDD '18)*. ACM, New York, NY, USA, 993–1001. https://doi.org/10.1145/3219819.3219820

[33] Zhi-Hua Zhou and Xu-Ying Liu. 2006. Training cost-sensitive neural networks with methods addressing the class imbalance problem. *IEEE Transactions on Knowledge and Data Engineering* 18, 1 (Jan 2006), 63–77. https://doi.org/10.1109/TKDE.2006.17

[34] Xiangxin Zhu, Dragomir Anguelov, and Deva Ramanan. 2014. Capturing long-tail distributions of object subcategories. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 915–922.