

Cost-Effective Active Learning for Deep Image Classification

Keze Wang, Dongyu Zhang, Ya Li, Ruimao Zhang, and Liang Lin, *Senior Member, IEEE*

Abstract—Recent successes in learning-based image classification, however, heavily rely on the large number of annotated training samples, which may require considerable human effort. In this paper, we propose a novel active learning (AL) framework, which is capable of building a competitive classifier with optimal feature representation via a limited amount of labeled training instances in an incremental learning manner. Our approach advances the existing AL methods in two aspects. First, we incorporate deep convolutional neural networks into AL. Through the properly designed framework, the feature representation and the classifier can be simultaneously updated with progressively annotated informative samples. Second, we present a cost-effective sample selection strategy to improve the classification performance with less manual annotations. Unlike traditional methods focusing on only the uncertain samples of low prediction confidence, we especially discover the large amount of high-confidence samples from the unlabeled set for feature learning. Specifically, these high-confidence samples are automatically selected and iteratively assigned pseudolabels. We thus call our framework cost-effective AL (CEAL) standing for the two advantages. Extensive experiments demonstrate that the proposed CEAL framework can achieve promising results on two challenging image classification data sets, i.e., face recognition on the cross-age celebrity face recognition data set database and object categorization on Caltech-256.

Index Terms—Active learning (AL), deep neural nets, image classification, incremental learning.

I. INTRODUCTION

AIMING at improving the existing models by incrementally selecting and annotating the most informative unlabeled samples, active learning (AL) has been well studied in the past few decades [3]–[12], and applied to various kinds

of vision tasks, such as image/video categorization [13]–[17], text/Web classification [18]–[20], and image/video retrieval [21], [22]. In the AL methods [3]–[5], the classifier/model is first initialized with a relatively small set of labeled training samples. Then it is continuously boosted by selecting and pushing some of the most informative samples to user for annotation. Although the existing AL approaches [10], [11] have demonstrated impressive results on image classification, their classifiers/models are trained with hand-craft features (e.g., HoG and SIFT) on small-scale visual data sets. The effectiveness of AL on more challenging image classification tasks has not been studied well.

Recently, incredible progress on visual recognition tasks has been made by deep learning approaches [23], [24]. With sufficient labeled data [25], deep convolutional neural networks (CNNs) [23], [26] are trained to directly learn features from raw pixels, which have achieved the state-of-the-art performance for image classification. However, in many real applications of large-scale image classification, the labeled data are not enough, since the tedious manual labeling process requires a lot of time and labor. Thus, it has a great practical significance to develop a framework by combining CNNs and AL, which can jointly learn features and classifiers/models from unlabeled training data with minimal human annotations. However, incorporating CNNs into AL framework is not straightforward for real image classification tasks. This is due to the following two issues.

- 1) The labeled training samples given by current AL approaches are insufficient for CNNs, as the majority of unlabeled samples are usually ignored in AL. AL usually selects only a few of the most informative samples (e.g., samples with quite low prediction confidence) in each learning step and frequently solicit user labeling. Thus, it is difficult to obtain proper feature representations by fine-tuning CNNs with these minority of informative samples.
- 2) The process pipelines of AL and CNNs are inconsistent with each other. Most of AL methods pay close attention to model/classifier training. Their strategies to select the most informative samples are heavily dependent on the assumption that the feature representation is fixed. However, the feature learning and classifier training are jointly optimized in CNNs. Because of this inconsistency, simply fine-tuning CNNs in the traditional AL framework may face the divergence problem.

Inspired by the insights and lessons from a significant amount of previous works as well as the recently

Manuscript received June 26, 2015; revised January 5, 2016 and April 25, 2016; accepted July 1, 2016. Date of publication July 11, 2016; date of current version December 13, 2017. This work was supported in part by the National Natural Science Foundation of China under Grant 61622214, in part by the State Key Development Program under Grant 2016YFB1001000, in part sponsored by CCF-Tencent Open Fund, in part by the Special Program through the Applied Research on Super Computation of the Natural Science Foundation of China–Guangdong Joint Fund (the second phase), and in part by NVIDIA Corporation through the Tesla K40 GPU. This paper was recommended by Associate Editor E. Cetin. (*Corresponding author: Dongyu Zhang.*)

K. Wang, D. Zhang, R. Zhang, and L. Lin are with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou 510006, China, and also with Collaborative Innovation Center of High Performance Computing, National University of Defense Technology, Changsha 410073, China. (e-mail: kezewang@gmail.com; cszhangdy@163.com; r.m.zhang1989@gmail.com; linliang@ieee.org).

Y. Li is with Guangzhou University, Guangzhou 510182, China (e-mail: liya@gzhu.edu.cn).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCSVT.2016.2589879

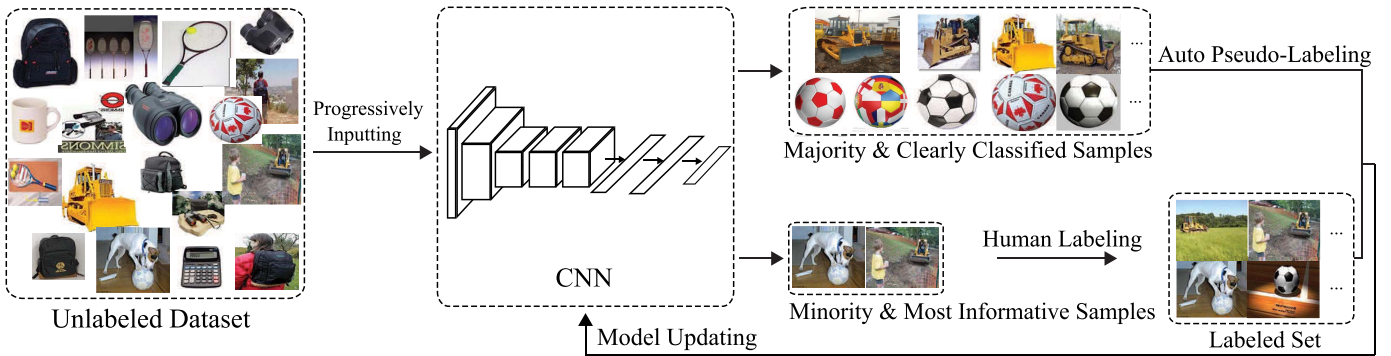


Fig. 1. Illustration of our proposed CEAL framework. Our proposed CEAL progressively feeds the samples from the unlabeled data set into the CNN. Then both of the clearly classified samples and most informative samples selection criteria are applied on the classifier output of the CNN. After adding user-annotated minority of uncertain samples into the labeled set and pseudolabeling the majority of certain samples, the model (feature representation and classifier of the CNN) is further updated.

proposed technique, i.e., self-paced learning [27]–[30], we address above-mentioned issues by cost effectively combining the CNN and AL via a complementary sample selection. In particular, we propose a novel AL framework called cost-effective AL (CEAL), which is enabled to fine-tune the CNN with sufficient unlabeled training data and overcomes the inconsistency between the AL and CNN.

Different from the existing AL approaches that consider only the most informative and representative samples, our CEAL proposes to automatically select and pseudoannotate unlabeled samples. As Fig. 1 illustrates, our proposed CEAL progressively feeds the samples from the unlabeled data set into the CNN and selects two kinds of samples for fine-tuning according to the output of CNN’s classifiers. One kind is the minority of samples with low prediction confidence, called most informative/uncertain samples. The predicted labels of samples are most uncertainty ones. For the selection of these uncertain samples, the proposed CEAL considers three common AL methods: least confidence (LC) [31], margin sampling (MS) [32], and entropy (EN) [33]. The selected samples are added into the labeled set after active user labeling. The other kind is the majority of samples with high prediction confidence, called high-confidence samples. The predicted labels of samples are most certainty ones. For these certain kinds of samples, the proposed CEAL automatically assigns pseudolabels with no human labor cost. As one can see, these two kinds of samples are complementary to each other for representing different confidence levels of the current model on the unlabeled data set. In the model updating stage, all the samples in the labeled set and currently pseudolabeled high-confidence samples are exploited to fine-tune the CNN.

The proposed CEAL advances in employing these two complementary kinds of samples to incrementally improve the model’s classifier training and feature learning: the minority of informative kind contributes to train more powerful classifiers, while the majority of high confidence kind conduces to learn more discriminative feature representations. On one hand, although the number is small, most uncertainty unlabeled samples usually have a great potential impact on the classifiers. Selecting and annotating them into training can lead to a better decision boundary of the classifiers. On the other hand, though unable to significantly improve the performance of

classifiers, the high-confidence unlabeled samples are close to the labeled samples in the CNN’s feature space. Thus, pseudolabeling these majority of high-confidence samples for training is a reasonable data augmentation way for the CNN to learn robust features. In particular, the number of the high-confidence samples is actually much larger than that of most uncertainty ones. With the obtained robust feature representation, the inconsistency between the AL and CNN can be overcome.

For the problem of keep the model stable in the training stage, many works [34], [35] are proposed in recent years inspired by the learning process of humans that gradually include samples into training from easy to complex. Through this way, the training samples for further iterations are gradually determined by the model itself based on what it has already learned [30]. In other words, the model can gradually select the high-confidence samples as pseudolabeled ones along with the training process. The advantages of these related studies motivate us to incrementally select unlabeled samples in an easy-to-hard manner to make pseudolabeling process reliable. Specifically, considering that the classification model is usually not reliable enough in the initial iterations, we employ high-confidence threshold to define clearly classified samples and assign them pseudolabels. When the performance of the classification model improves, the threshold correspondingly decreases.

The main contribution of this paper is threefold. First, to the best of our knowledge, our work is the first one addressing the deep image classification problems in conjunction with AL framework and CNN training. Our framework can be easily extended to other similar visual recognition tasks. Second, this paper also advances the AL development by introducing a cost-effective strategy to automatically select and annotate the high-confidence samples, which improves the traditional samples selection strategies. Third, experiments on challenging cross-age celebrity face recognition data set (CACD) [1] and Caltech 256 [2] data sets show that our approach outperforms other methods not only in the classification accuracy but also in the reduction of human annotation.

The rest of this paper is organized as follows. Section II presents a brief review of related work. Section III discusses the component of our framework and the corresponding

learning algorithm. Section IV presents the experiments results with deep empirical analysis. Section V concludes this paper.

II. RELATED WORK

The key idea of the AL is that a learning algorithm should achieve higher accuracy with a fewer labeled training samples, if it is allowed to choose the ones from which it learns [31]. In this way, the instance selection scheme is becoming extremely important. One of the most common strategy is the uncertainty-based selection [12], [18], which measures the uncertainties of novel unlabeled samples from the predictions of previous classifiers. Lewis [12] proposed to extract the sample, which has the largest EN on the conditional distribution over predicted labels, as the most uncertain instance. The support vector machine (SVM)-based method [18] determined the uncertain samples based on the relative distance between the candidate samples and the decision boundary. Some earlier works [19], [38] also determined the sample uncertainty referring to a committee of classifiers (i.e., examining the disagreement among class labels assigned by a set of classifiers). Such a theoretically motivated framework is called *query-by-committee* in literature [31]. All the above-mentioned uncertainty-based methods usually ignore the majority of certain unlabeled samples and thus are sensitive to outliers. The latter methods have taken the information density measure into account and exploited the information of unlabeled data when selecting samples. These approaches usually sequentially select the informative samples relying on the probability estimation [6], [37] or prior information [8] to minimize the generalization error of the trained classifier over the unlabeled data. For example, Joshi *et al.* [6] considered the uncertainty sampling method based on the probability estimation of class membership for all the instances in the selection pool, and such a method can be effective to handle the multiclass case. In [8], some context constraints are introduced as the priori to guide users to tag the face images more efficiently. At the same time, a series of works [7], [24] is proposed to take the samples to maximize the increase of mutual information between the candidate instance and the remaining unlabeled instances under the Gaussian process framework. Li and Guo [10] presented a novel adaptive AL approach that combines an information density measure and a most uncertainty measure together to label critical instances for image classifications. Moreover, the diversity of the selected instance over the certain category has been taken into consideration in [4] as well. Such a work is also the pioneer study expanding the SVM-based AL from the *single mode* to *batch mode*. Recently, Elhamifar *et al.* [11] further integrated the uncertainty and diversity measurement into a unified *batch mode* framework via convex programming for unlabeled sample selection. Such an approach is more feasible to conjunction with any type of classifiers, but not limited in max-margin based ones. It is obvious that all the above-mentioned AL methods consider only those low-confidence samples (e.g., uncertain and diverse samples), but losing the sight of a large majority of high-confidence samples. We hold that due to the majority and consistency, these high-confidence

samples will also be beneficial to improve the accuracy and keep the stability of classifiers. Even more, we shall demonstrate that considering these high-confidence samples can also reduce the user effort of annotation effectively.

III. COST-EFFECTIVE ACTIVE LEARNING

In this section, we develop an efficient algorithm for the proposed CEAL framework. Our objective is to apply our CEAL framework to deep image classification tasks by progressively selecting complementary samples for model updating. Suppose we have a data set of m categories and n samples denoted by $D = \{x_i\}_{i=1}^n$. We denote the currently annotated samples of D by D^L , while the unlabeled ones by D^U . The label of x_i is denoted by $y_i = j$, $j \in \{1, \dots, m\}$, i.e., x_i belongs to the j th category. We should give two necessary remarks on our problem settings. One is that in our investigated image classification problems, almost all data are unlabeled, i.e., most of the $\{y_i\}$ values of D are unknown and needed to be completed in the learning process. The other remark is that D^U might possibly been input into the system in an incremental way. This means that data scale might be consistently growing.

Thanks to handling both manually annotated and automatically pseudolabeled samples together, our proposed CEAL model can progressively fit the consistently growing unlabeled data in such a holistic manner. The CEAL for deep image classification is formulated as follows:

$$\min_{\{\mathcal{W}, y_i, i \in D^U\}} -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \mathbf{1}\{y_i = j\} \log p(y_i = j | x_i; \mathcal{W}) \quad (1)$$

where $\mathbf{1}\{\cdot\}$ is the indicator function, so that $\mathbf{1}\{\text{a true statement}\} = 1$ and $\mathbf{1}\{\text{a false statement}\} = 0$, and \mathcal{W} denotes the network parameters of the CNN. $p(y_i = j | x_i; \mathcal{W})$ denotes the softmax output of the CNN for the j th category, which represents the probability of the sample x_i belonging to the j th classifiers.

The alternative search strategy is readily employed to optimize (1). Specifically, the algorithm is designed by alternatively updating the pseudolabeled sample $y_i \in D^U$ and the network parameters \mathcal{W} . In the following, we introduce the details of the optimization steps and give their physical interpretations. The practical implementation of the CEAL will also be discussed in the end.

A. Initialization

Before the experiment starts, the labeled samples D^L is empty. For each class, we randomly select a few training samples from D^U and manually annotate them as the starting point to initialize the CNN parameters \mathcal{W} .

B. Complementary Sample Selection

Fixing the CNN parameters \mathcal{W} , we first rank all unlabeled samples according to the common AL criteria and then manually annotate those most uncertain samples and add them into D^L . For those most certain ones, we assign pseudolabels and denote them by D^H .

1) *Informative Sample Annotating*: Our CEAL can use in conjunction with any type of common actively learning criteria, e.g., LC [31], MS [32], and EN [33] to select K most informative/uncertain samples left in D^U . The selection criteria are based on $p(y_i = j|x_i; \mathcal{W})$, which denotes the probability of x_i belonging to the j th class. Specifically, the three selection criteria are defined as follows.

- 1) *LC*: Rank all the unlabeled samples in an ascending order according to the lc_i value. lc_i is defined as

$$lc_i = \max_j p(y_i = j|x_i; \mathcal{W}). \quad (2)$$

If the probability of the most probable class for a sample is low, then the classifier is uncertain about the sample.

- 2) *MS*: Rank all the unlabeled samples in an ascending order according to the ms_i value. ms_i is defined as

$$ms_i = p(y_i = j_1|x_i; \mathcal{W}) - p(y_i = j_2|x_i; \mathcal{W}) \quad (3)$$

where j_1 and j_2 represent the first and second most probable class labels predicted by the classifiers, respectively. The smaller of the margin means the classifier is more uncertain about the sample.

- 3) *EN*: Rank all the unlabeled samples in a descending order according to their en_i value. en_i is defined as

$$en_i = - \sum_{j=1}^m p(y_i = j|x_i; \mathcal{W}) \log p(y_i = j|x_i; \mathcal{W}). \quad (4)$$

This method takes all class label probabilities into consideration to measure the uncertainty. The higher EN value, the more uncertain is the sample.

2) *High-Confidence Sample Pseudolabeling*: We select the high-confidence samples from D^U , whose EN is smaller than the threshold δ . Then we assign clearly predicted pseudolabels to them. The pseudolabel y_i is defined as

$$y_i = \begin{cases} j^*, & en_i < \delta \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

where $y_i = 1$ denotes that x_i is regarded as high-confidence sample. The selected samples are denoted by D^H . Note that compared with classification probability $p(y_i = j^*|x_i; \mathcal{W})$ for the j^* th category, the employed EN en_i holistically considers the classification probability of the other categories, i.e., the selected sample should be clearly classified with high confidence. The threshold δ is set to a large value to guarantee a high reliability of assigning a pseudolabel.

C. CNN Fine-Tuning

Fixing the labels of self-labeled high-confidence samples D^H and manually annotated ones D^L by active user, (1) can be simplified as

$$\min_{\mathcal{W}} - \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m \mathbf{1}\{y_i = j\} \log p(y_i = j|x_i; \mathcal{W}) \quad (6)$$

where N denotes the number of samples in $D^H \cup D^L$. We employ the standard back propagation to update the CNN's parameters \mathcal{W} . Specifically, let \mathcal{L} denote the loss function of (6), then the partial derive of the network parameter \mathcal{W} according to (6) is

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathcal{W}} &= \frac{-\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m \mathbf{1}\{y_i = j\} \log p(y_i = j|x_i; \mathcal{W})}{\partial \mathcal{W}} \\ &= -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m \mathbf{1}\{y_i = j\} \frac{\partial \log p(y_i = j|x_i; \mathcal{W})}{\partial \mathcal{W}} \\ &= -\frac{1}{N} \sum_{i=1}^N (\mathbf{1}\{y_i = j\} - p(y_i = j|x_i; \mathcal{W})) \frac{\partial z_j(x_i; \mathcal{W})}{\partial \mathcal{W}} \end{aligned} \quad (7)$$

where $\{z_j(x_i; \mathcal{W})\}_{j=1}^m$ denotes the activation for the i th sample of the last layer of CNN model before feeding into the softmax classifier, which is defined as

$$p(y_i = j|x_i; \mathcal{W}) = \frac{e^{z_j(x_i; \mathcal{W})}}{\sum_{t=1}^m e^{z_t(x_i; \mathcal{W})}}. \quad (8)$$

After fine-tuning, we put the high-confidence samples D^H back to D^U and erase their pseudolabel.

D. Threshold Updating

As the incremental learning process goes on, the classification capability of classifier improves and more high-confidence samples are selected, which may result in the decrease of incorrect autoannotation. In order to guarantee the reliability of high-confidence sample selection, at the end of each iteration t , we update the high-confidence sample selection threshold by setting

$$\delta = \begin{cases} \delta_0, & t = 0 \\ \delta - dr * t, & t > 0 \end{cases} \quad (9)$$

where δ_0 is the initial threshold and dr controls the threshold decay rate.

The entire algorithm can be then summarized into Algorithm 1. It is easy to see that this alternative optimizing strategy finely accords with the pipeline of the proposed CEAL framework.

IV. EXPERIMENTS

A. Data Sets and Experimental Settings

1) *Data Set Description*: In this section, we evaluate our CEAL framework on two public challenging benchmarks, i.e., CACD [1] and the Caltech-256 object categorization [2] data set (see Fig. 2). CACD is a large-scale and challenging data set for face identification and retrieval problems. It contains more than 160 000 images of 2000 celebrities, which are varying in age, pose, illumination, and occlusion. Since not all of the images are annotated, we adopt a subset of 580 individuals from the whole data set in our experiments, in which 200 individuals are originally annotated and 380 persons are extra annotated by us. Especially, 6336 images of 80 individuals are utilized for pretraining the network and the remaining



Fig. 2. Demonstration of the effectiveness of our proposed heuristic deep AL framework on face recognition and object categorization. First and second lines: sample images from the Caltech-256 [2] data set. Last line: samples images from CACD [1].

Algorithm 1 Learning Algorithm of CEAL

Input:

Unlabeled samples D^U , initially labeled samples D^L , uncertain samples selection size K , high-confidence samples selection threshold δ , threshold decay rate dr , maximum iteration number T , fine-tuning interval t .

Output:

CNN parameters \mathcal{W} .

- 1: Initialize \mathcal{W} with D^L .
- 2: **while** not reach maximum iteration T **do**
- 3: Add K uncertainty samples into D^L based on Eq. (2) or (3) or (4),
- 4: Obtain high confidence samples D^H based on Eq. (5)
- 5: In every t iterations:
 - Update \mathcal{W} via fine-tuning according to Eq. (6) with $D^H \cup D^L$
 - Update δ according to Eq. (9)
- 6: **end while**
- 7: **return** \mathcal{W}

500 persons are used to perform the experiments. Caltech-256 is a challenging object categories data set. It contains a total of 30 607 images of 256 categories collected from the Internet.

2) *Experimental Setting*: For CACD, we utilize the method proposed in [38] to detect the facial points and align the faces based on the eye locations. We resize all the faces into 200×150 and then we set the parameters: $\delta_0 = 0.05$, $dr = 0.0033$, and $K = 2000$. For Caltech-256, we resize all the images to 256×256 and we set $\delta_0 = 0.005$, $dr = 0.00033$, and $K = 1000$. Following the settings in the existing AL method [11], we randomly select 80% images of each class to form the unlabeled training set, and the rest are as the testing set in our experiments. Among the unlabeled training set, we randomly select 10% samples of each class to initialize the network and the rest are for incremental learning process. To get rid of the influence of randomness, we average five times execution results as the final result.

We use different network architectures for CACD [1] and Caltech-256 [2] data sets because the difference between face and object is relatively large. Table I shows the overall network

TABLE I

DETAILED CONFIGURATION OF THE CNN ARCHITECTURE USED IN CACD [1]. IT TAKES THE $200 \times 150 \times 3$ IMAGES AS INPUT AND GENERATES THE 500-WAY SOFTMAX OUTPUT FOR CLASSES PREDICTION. THE ReLU [39] ACTIVATION FUNCTION IS NOT SHOWN FOR BREVITY

layer type	kernel size/stride	output size
convolution	$5 \times 5/2$	$98 \times 73 \times 32$
max pool	$3 \times 3/2$	$48 \times 36 \times 32$
LRN		$48 \times 36 \times 32$
convolution(padding2)	$5 \times 5/1$	$48 \times 36 \times 64$
max pool	$3 \times 3/2$	$23 \times 17 \times 64$
LRN		$23 \times 17 \times 64$
convolution(padding1)	$3 \times 3/1$	$23 \times 17 \times 96$
fc(dropout50%)		$1 \times 1 \times 1536$
fc(dropout50%)		$1 \times 1 \times 1536$
softmax		$1 \times 1 \times 500$

TABLE II

DETAILED CONFIGURATION OF THE CNN ARCHITECTURE USED IN CALTECH-256 [2]. IT TAKES THE $256 \times 256 \times 3$ IMAGES AS INPUT, WHICH WILL BE RANDOMLY CROPPED INTO 227×227 DURING THE TRAINING, AND GENERATES THE 256-WAY SOFTMAX OUTPUT FOR CLASS PREDICTION. THE ReLU ACTIVATION FUNCTION IS NOT SHOWN FOR BREVITY

layer type	kernel size/stride	output size
convolution	$11 \times 11/4$	$55 \times 55 \times 96$
max pool	$3 \times 3/2$	$27 \times 27 \times 96$
LRN		$27 \times 27 \times 96$
convolution(padding2)	$5 \times 5/1$	$27 \times 27 \times 256$
max pool	$3 \times 3/2$	$13 \times 13 \times 256$
LRN		$13 \times 13 \times 256$
convolution(padding1)	$3 \times 3/1$	$13 \times 13 \times 384$
convolution(padding1)	$3 \times 3/1$	$13 \times 13 \times 384$
convolution(padding1)	$3 \times 3/1$	$13 \times 13 \times 256$
max pool	$3 \times 3/2$	$6 \times 6 \times 256$
fc(dropout50%)		$1 \times 1 \times 4096$
fc(dropout50%)		$1 \times 1 \times 4096$
softmax		$1 \times 1 \times 256$

architecture for CACD experiments, and Table II shows the overall network architecture for Caltech-256 experiments. We use Alexnet [23] as the network architecture for Caltech-256 and using the ImageNet ILSVRC data set [40] pretrained model as the starting point following the setting of [41]. Then we keep all layers fixed and just modify the last layer to be the 256-way softmax classifier to perform the Caltech-256 experiments. We employ Caffe [42] for CNN implementation.

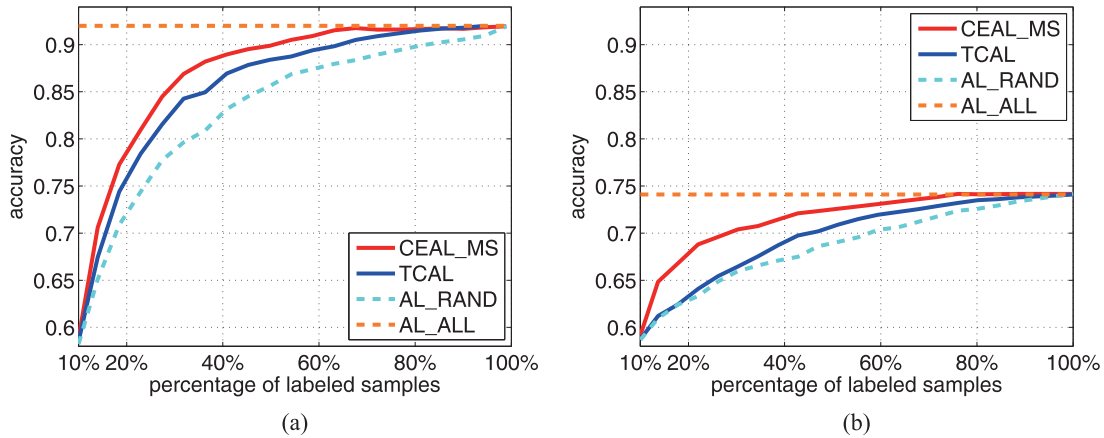


Fig. 3. Classification accuracy under different percentages of annotated samples of the whole training set on the (a) CACD and (b) Caltech-256 data sets. Our proposed method CEAL_MS performs consistently better than the compared TCAL and AL_RAND.

For CACD, we set the learning rates of all the layers as 0.01. For Caltech-256, we set the learning rates of all the layers as 0.001 except for the softmax layer, which is set to 0.01. All the experiments are conducted on a common desktop PC with an intel 3.8-GHz CPU and a Titan X GPU. Average 17 h are needed to finish training on the CACD data set with 44 708 images.

3) *Comparison Methods*: To demonstrate that our proposed CEAL framework can improve the classification performance with less labeled data, we compare CEAL with new state-of-the-art AL [triple criteria AL (TCAL)] and baseline methods (AL_ALL and AL_RAND).

- 1) *AL_ALL*: We manually label all the training samples and use them to train the CNN. This method can be regarded as the upper bound (best performance that CNN can reach with all labeled training samples).
- 2) *AL_RAND*: During the training process, we randomly select samples to be annotated to fine-tune the CNN. This method discards all AL techniques and can be considered as the lower bound.
- 3) *TCAL [3]*: TCAL is a comprehensive AL approach and is well designed to jointly evaluate sample selection criteria (uncertainty, diversity and density), and has overcome the state-of-the-art methods with much less annotations. TCAL represents those methods who intend to mine minority of informative samples to improve the performance. Thus, we regard it as a relevant competitor.

Implementation Details: The compared methods share the same CNN architecture with our CEAL on the both data sets. The only difference in the sample selection criteria. For the BaseLine method, we select all training samples to fine-tune the CNN, i.e., all labels are used. For TCAL, we follow the pipeline of [3] by training an SVM classifier and then applying the uncertainty, diversity and density criteria to select the most informative samples. Specifically, the uncertainty of samples is assessed according to the MS strategy. The diversity is calculated by clustering the most uncertain samples via k-means with histogram intersection kernel. The density of one sample is measured by calculating the average distance with other samples within a cluster it belonged to. For each cluster,

the highest density (i.e., the smallest average distance) sample is selected as the most informative sample. For CACD, we cluster 2000 most uncertain samples and select 500 most informative samples according to the above-mentioned diversity and density. For Caltech-256, we select 250 most informative samples from 1000 most uncertain samples. To make a fair comparison, samples selected in each iteration by the TCAL are also used to fine-tune the CNN to learn the optimal feature representation as AL_RAND. Once optimal feature learned, the SVM classifier of TCAL is further updated.

B. Comparison Results and Empirical Analysis

1) *Comparison Results*: To demonstrate the effectiveness of our proposed framework, we also apply the MS criterion to measure the uncertainty of samples and denote this method by CEAL_MS. Fig. 3 illustrates the accuracy-percentage of annotated samples curve of AL_RAND, AL_ALL, TCAL, and the proposed CEAL_MS on both CACD and Caltech-256 data sets. This curve demonstrates the classification accuracy under different percentages of annotated samples of the whole training set.

As illustrated in Fig. 3, Table III(a) and (b), our proposed CEAL framework overcomes the compared method from the aspects of the recognition accuracy and user annotation amount. From the aspect of recognition accuracy, given the same percentage of annotated samples, our CEAL_MS outperforms the compared method in a clear margin, especially when the percentage of annotated samples is low. From the aspect of the user annotation amount, to achieve 91.5% recognition accuracy on the CACD data set, AL_RAND and TCAL require 99% and 81% labeled training samples, respectively. CEAL_MS needs only 63% labeled samples and reduces around 36% and 18% user annotations, compared with AL_RAND and TCAL. To achieve the 73.8% accuracy on the caltech-256 data set, AL_RAND and TCAL require 97% and 93% labeled samples, respectively. CEAL_MS needs only 78% labeled samples and reduces around 19% and 15% user annotations, compared with AL_RAND and TCAL. This justifies that our proposed CEAL framework can effectively reduce the need of labeled samples.

TABLE III
CLASS ACCURACY PER SOME SPECIFIC AL ITERATIONS ON THE (a) CACD AND (b) CALTECH-256 DATA SETS

(a)

Training iteration	0	2	4	6	8	10	12	14	16	18	20
Percentage of labeled samples	0.1	0.18	0.27	0.36	0.45	0.54	0.63	0.72	0.81	0.90	1
CEAL_MS	57.4%	77.3%	84.5%	88.2%	89.5%	90.5%	91.5%	91.6%	91.7%	91.7%	92.0%
TCAL	57.4%	74.4%	81.6%	85.0%	87.9%	88.8%	89.8%	90.9%	91.5%	91.8%	91.9%
AL_RAND	57.4%	70.9%	77.7%	80.9%	84.5%	86.9%	88.0%	89.0%	89.9%	90.6%	92.0%
AL_ALL	-	-	-	-	-	-	-	-	-	-	92.0%

(b)

Training iteration	0	2	4	6	8	10	12	14	16	18	20	22
Percentage of labeled samples	0.10	0.18	0.26	0.34	0.43	0.51	0.59	0.68	0.76	0.84	0.93	1
CEAL_MS	58.4%	64.8%	68.8%	70.4%	71.4%	72.3%	72.8%	73.3%	73.8%	74.2%	74.2%	74.2%
TCAL	58.4%	62.4%	65.4%	67.5%	69.8%	70.9%	71.9%	72.5%	73.2%	73.6%	73.9%	74.2%
AL_RAND	58.4%	62.4%	64.8%	66.6%	67.5%	69.1%	70.4%	71.2%	72.4%	72.9%	73.6%	74.2%
AL_ALL	-	-	-	-	-	-	-	-	-	-	-	74.2%

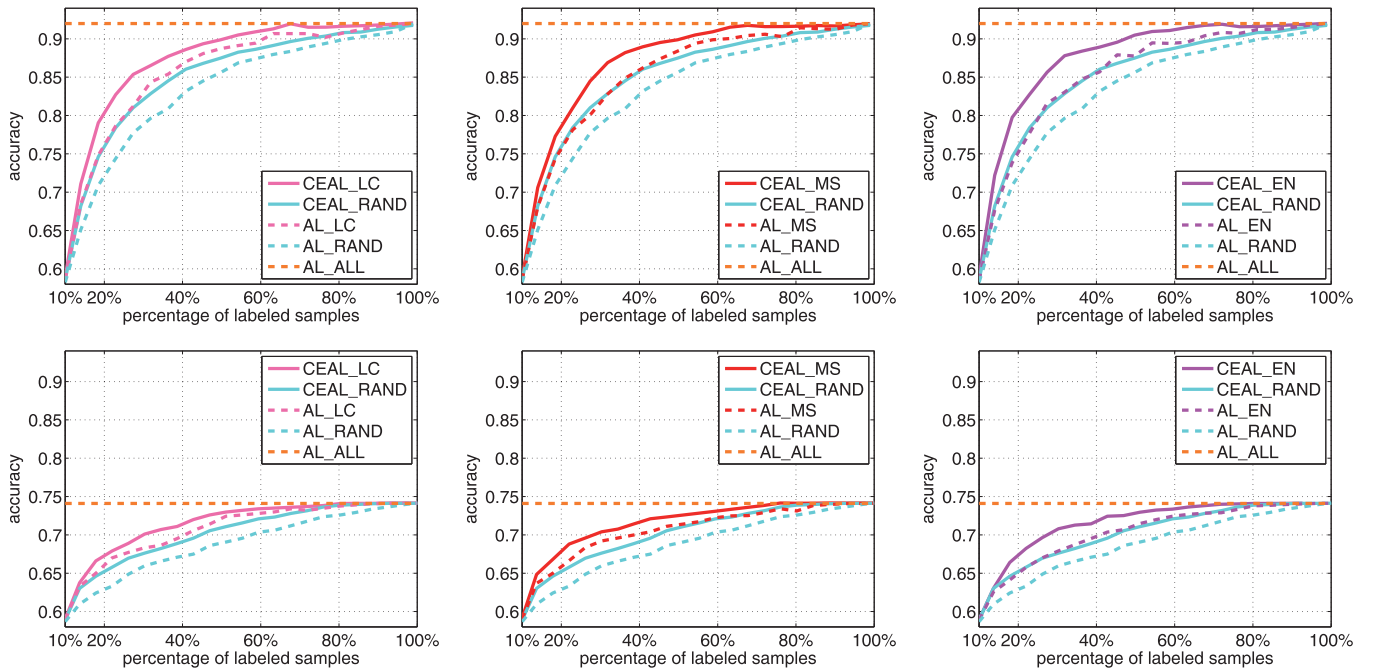


Fig. 4. Extensive study for different informative sample selection criteria on CACD (the first row) and Caltech-256 (the second row) data sets. These criteria include LC (the first column), MS (the second column), and EN (the third column). One can observe that our CEAL framework works consistently well on the common information sample selection criteria.

From the above results, one can see that our proposed frame CEAL performs consistently better than the state-of-the-art method TCAL in both recognition accuracy and user annotation amount through fair comparisons. This is due to that TCAL only mines minority of informative samples and is not able to provide sufficient training data for feature learning under the deep image classification scenario. Hence, our CEAL has a competitive advantage in deep image classification task. To clearly analyze our CEAL and justify the effectiveness of its component, we have conducted the several experiments and discussed in the following sections.

2) *Component Analysis*: To justify that the proposed CEAL can work consistently well on the common informative sample selection criteria, we implement three variants of CEAL according to LC, MS, and EN to assess uncertain samples. These three variants are denoted by CEAL_LC, CEAL_MS,

and CEAL_EN. Meanwhile, to show the raw performance of these criteria, we discard the cost-effective high-confidence sample selection of the above-mentioned variants and denoted the discarded versions by AL_LC, AL_MS, and AL_EN. To clarify the contribution of our pseudolabeling majority of high-confidence samples strategy, we further introduce this strategy into the AL_RANDOM and denote this variant by CEAL_RANDOM. Since AL_RANDOM means randomly select samples to be annotated, CEAL_RANDOM reflects the original contribution of the pseudolabeled majority of high-confidence samples strategy, i.e., CEAL_RANDOM denotes the method that uses only the pseudolabeled majority of samples.

Fig. 4 illustrates the results of these variants on the data sets CACD (the first row) and Caltech-256 (the second row). The results demonstrate that giving the same percentage of labeled samples and compared with AL_RANDOM, CEAL_RANDOM,

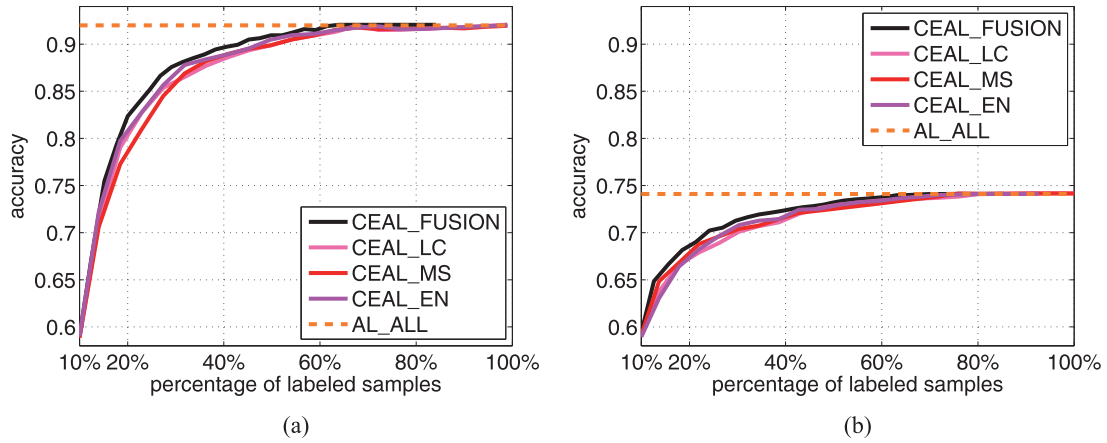


Fig. 5. Comparison between different informative sample selection criteria and their fusion (CEAL_FUSION) on (a) CACD and (b) Caltech-256 data sets.

simply exploiting pseudolabeled majority samples, obtains similar performance gain as AL_LC, AL_MS, and AL_EN, which employs the informative sample selection criterion. This justifies that our proposed pseudolabeled majority of samples strategy is effective as some common informative sample selection criteria. Moreover, as one can see that in Fig. 4, CEAL_LC, CEAL_MS, and CEAL_EN all consistently outperform the pure pseudolabeling samples version CEAL_RANDOM and their excluding pseudolabeled samples versions AL_LC, AL_MS, and AL_EN in a clear margin on both the CACD and Caltech-256 data sets, respectively. This validates that our proposed pseudolabeled majority of samples strategy is complementary to the common informative sample selection criteria and can further significantly improve the recognition performance.

To analyze the choice of informative sample selection criteria, we have made a comparison among the three above-mentioned criteria. We also make an attempt to simply combine them together. Specifically, in each iteration, we select top $K/2$ samples according to each criterion, respectively. Then we remove repeated ones (i.e., some samples may be selected by the three criteria at the same time) from the obtained $3K/2$ samples. After removing the repeated samples, we randomly select K samples from them to require user annotations. We denote this method by CEAL_FUSION.

Fig. 5 illustrates that CEAL_LC, CEAL_MS, and CEAL_EN have a similar performance, while CEAL_FUSION performs better. This demonstrates that the informative sample selection criterion still plays an important role in improving the recognition accuracy. Though being a minority, the informative samples have a great potential impact on the classifier.

C. Reliability of CEAL

From the above experiments, we know that the performance of our framework is better than those of other methods, which shows the superiority of introducing the majority of pseudolabeled samples. But how does the accuracy of assigning the pseudo-label to those high-confidence samples? In order to demonstrate the reliability of our proposed CEAL framework, we also evaluate the average error in selecting

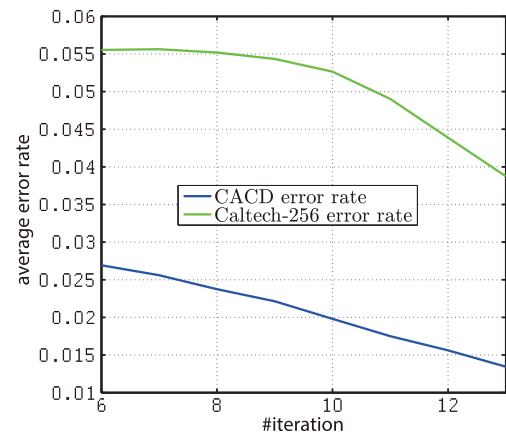


Fig. 6. Average error rate of the pseudolabels of high-confidence samples assigned by the heuristic strategy on the CACD and Caltech-256 data sets experiments. The vertical axes represent the average error rate and the horizontal axes represent the learning iteration. Our proposed CEAL framework can assign reliable pseudolabels to the unlabeled samples under acceptable average error rate.

high-confidence samples. Fig. 6 shows the error rate of assigning pseudolabel along with the learning iteration. As one can see, the average error rate is quite low (say less than 3% on the CACD data set and less than 5.5% on the Caltech-256 data set) even at early iterations. Hence, our proposed CEAL framework can assign reliable pseudolabels to the unlabeled samples under acceptable average error rate along with the learning iteration.

D. Sensitivity of High-Confidence Threshold

Since the training phase of deep CNNs is time consuming, it is not affordable to employ a try and error approach to set the threshold for defining high-confidence samples. We further analyze the sensitivity of the threshold parameters δ (threshold) and dr (threshold decay rate) on our system performance on the CACD data set using CEAL_EN. While analyzing the sensitivity of the parameter δ on our system, we fix the decrease rate dr to 0.0033. We fix the threshold δ to 0.05 when analyzing the sensitivity of dr . The results of the sensitivity analysis of δ (range 0.045 to 0.1) are

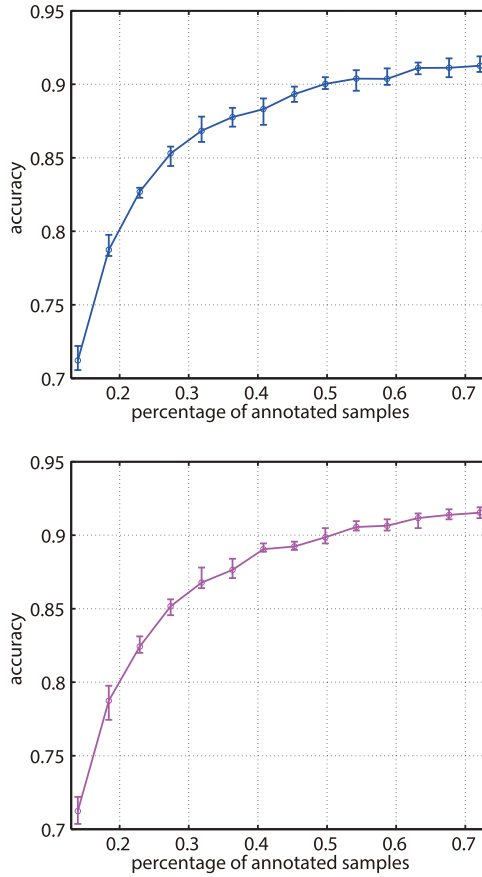


Fig. 7. Sensitivity analysis of heuristic threshold δ (top) and decay rate dr (bottom). One can observe that these parameters do not substantially affect the overall system performance.

shown in the top of Fig. 7, while the sensitivity analysis of dr (range 0.001 to 0.0035) is shown in the bottom of Fig. 7. Note that the test range of δ and dr is set to ensure the majority of high confidence assumption of this paper. Though the range of $\{\delta, dr\}$ seems to be narrow from the value, it leads to a significant difference: about 10%–60% samples are pseudolabeled in high-confidence sample selection. The lower standard deviation of the accuracy in Fig. 7 proves that the choice of these parameters does not significantly affect the overall system performance.

V. CONCLUSION

In this paper, we propose a CEAL framework for deep image classification tasks, which employs a complementary sample selection strategy: progressively select the minority of most informative samples and pseudolabel the majority of high-confidence samples for model updating. In such a holistic manner, the minority of labeled samples benefit the decision boundary of classifier and the majority of pseudolabeled samples provide sufficient training data for robust feature learning. Extensive experiment results on two public challenging benchmarks justify the effectiveness of our proposed CEAL framework. In future work, we plan to apply our framework on more challenging large-scale object recognition tasks (e.g., 1000 categories in ImageNet). And we plan to incorporate more persons from the

CACD data set to evaluate our framework. Moreover, we plan to generalize our framework into other multilabel object recognition tasks (e.g., 20 categories in PASCAL VOC).

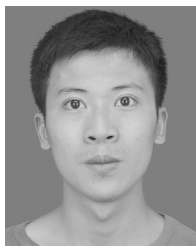
ACKNOWLEDGMENT

The authors would like to thank D. Liang and J. Xu for their preliminary contributions on this project.

REFERENCES

- [1] B.-C. Chen, C.-S. Chen, and W. H. Hsu, "Cross-age reference coding for age-invariant face recognition and retrieval," in *Proc. ECCV*, 2014, pp. 768–783.
- [2] G. Griffin, A. Holub, and P. Perona, "Caltech-256 object category dataset," California Inst. Technol., Pasadena, CA, USA, Tech. Rep. 7694, 2007.
- [3] B. Demir and L. Bruzzone, "A novel active learning method in relevance feedback for content-based remote sensing image retrieval," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 5, pp. 2323–2334, May 2015.
- [4] K. Brinker, "Incorporating diversity in active learning with support vector machines," in *Proc. ICML*, 2003, pp. 1–8.
- [5] B. Long, J. Bian, O. Chapelle, Y. Zhang, Y. Inagaki, and Y. Chang, "Active learning for ranking through expected loss optimization," *IEEE Trans. Knowl. Data Eng.*, vol. 27, no. 5, pp. 1180–1191, May 2015.
- [6] A. J. Joshi, F. Porikli, and N. Papanikolopoulos, "Multi-class active learning for image classification," in *Proc. CVPR*, Jun. 2009, pp. 2372–2379.
- [7] A. Kapoor, K. Grauman, R. Urtasun, and T. Darrell, "Active learning with Gaussian processes for object categorization," in *Proc. ICCV*, Oct. 2007, pp. 1–8.
- [8] A. Kapoor, G. Hua, A. Akbarzadeh, and S. Baker, "Which faces to tag: Adding prior constraints into active learning," in *Proc. ICCV*, Sep./Oct. 2009, pp. 1058–1065.
- [9] R. M. Castro and R. D. Nowak, "Minimax bounds for active learning," *IEEE Trans. Inf. Theory*, vol. 54, no. 5, pp. 2339–2353, May 2008.
- [10] X. Li and Y. Guo, "Adaptive active learning for image classification," in *Proc. CVPR*, Jun. 2013, pp. 859–866.
- [11] E. Elhamifar, G. Sapiro, A. Yang, and S. S. Sarsy, "A convex optimization framework for active learning," in *Proc. Int. Conf. Comput. Vis. (ICCV)*, Dec. 2013, pp. 209–216.
- [12] D. D. Lewis, "A sequential algorithm for training text classifiers: Corrigendum and additional data," *ACM SIGIR Forum*, vol. 29, no. 2, pp. 13–19, 1995.
- [13] X. Li and Y. Guo, "Multi-level adaptive active learning for scene classification," in *Proc. ECCV*, 2014, pp. 234–249.
- [14] B. Zhang, Y. Wang, and F. Chen, "Multilabel image classification via high-order label correlation driven active learning," *IEEE Trans. Image Process.*, vol. 23, no. 3, pp. 1430–1441, Mar. 2014.
- [15] F. Sun, M. Xu, and X. Jiang, "Robust multi-label image classification with semi-supervised learning and active learning," in *Proc. 21st Int. Conf. MultiMedia Modeling*, 2015, pp. 512–523.
- [16] S. C. H. Hoi, R. Jin, J. Zhu, and M. R. Lyu, "Batch mode active learning and its application to medical image classification," in *Proc. ICML*, 2006, pp. 417–424.
- [17] D. Tuia, F. Ratle, F. Pacifici, M. F. Kanevski, and W. J. Emery, "Correction to 'active learning methods for remote sensing image classification,'" *IEEE Trans. Geosci. Remote Sens.*, vol. 48, no. 6, p. 2767, Jun. 2010.
- [18] S. Tong and D. Koller, "Support vector machine active learning with applications to text classification," *J. Mach. Learn. Res.*, vol. 2, pp. 45–66, Mar. 2001.
- [19] A. McCallum and K. Nigam, "Employing EM and pool-based active learning for text classification," in *Proc. ICML*, 1998, pp. 350–358.
- [20] G. Schohn and D. Cohn, "Less is more: Active learning with support vector machines," in *Proc. ICML*, 2000, pp. 1–8.
- [21] S. Vijayanarasimhan and K. Grauman, "Large-scale live active learning: Training object detectors with crawled data and crowds," in *Proc. CVPR*, Jun. 2011, pp. 1449–1456.
- [22] A. G. Hauptmann, W.-H. Lin, R. Yan, J. Yang, and M.-Y. Chen, "Extreme video retrieval: Joint maximization of human and computer performance," in *Proc. 14th ACM Int. Conf. Multimedia*, 2006, pp. 385–394.

- [23] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. NIPS*, 2012, pp. 1097–1105.
- [24] D. Cireşan, U. Meier, and J. Schmidhuber, "Multi-column deep neural networks for image classification," in *Proc. CVPR*, Jun. 2012, pp. 3642–3649.
- [25] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. CVPR*, Jun. 2009, pp. 248–255.
- [26] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. ICLR*, 2015, pp. 1–14.
- [27] L. Jiang, D. Meng, Q. Zhao, S. Shan, and A. G. Hauptmann, "Self-paced curriculum learning," in *Proc. AAAI*, 2015, pp. 1–7.
- [28] Q. Zhao, D. Meng, L. Jiang, Q. Xie, Z. Xu, and A. G. Hauptmann, "Self-paced learning for matrix factorization," in *Proc. AAAI*, 2015, pp. 3196–3202.
- [29] L. Jiang, D. Meng, T. Mitamura, and A. G. Hauptmann, "Easy samples first: Self-paced reranking for zero-example multimedia search," in *Proc. 22nd ACM Int. Conf. Multimedia*, 2014, pp. 547–556.
- [30] L. Jiang, D. Meng, S.-I. Yu, Z. Lan, S. Shan, and A. G. Hauptmann, "Self-paced learning with diversity," in *Proc. NIPS*, 2014, pp. 2078–2086.
- [31] B. Settles, "Active learning literature survey," Comput. Sci. Dept., Univ. Wisconsin–Madison, Madison, WI, USA, Tech. Rep. 1648, 2009.
- [32] T. Scheffer, C. Decomain, and S. Wrobel, "Active hidden Markov models for information extraction," in *Proc. IDA*, 2001, pp. 309–318.
- [33] C. E. Shannon, "A mathematical theory of communication," *ACM SIGMOBILE Mobile Comput. Commun. Rev.*, vol. 5, no. 1, pp. 3–55, 2001.
- [34] M. P. Kumar, B. Packer, and D. Koller, "Self-paced learning for latent variable models," in *Proc. NIPS*, 2010, pp. 1189–1197.
- [35] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proc. ICML*, 2009, pp. 41–48.
- [36] Y. Freund, H. S. Seung, E. Shamir, and N. Tishby, "Selective sampling using the query by committee algorithm," *Mach. Learn.*, vol. 28, no. 2, pp. 133–168, 1997.
- [37] A. J. Joshi, F. Porikli, and N. P. Papanikolopoulos, "Scalable active learning for multiclass image classification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 11, pp. 2259–2273, Nov. 2012.
- [38] X. Xiong and F. De la Torre, "Supervised descent method and its applications to face alignment," in *Proc. CVPR*, Jun. 2013, pp. 532–539.
- [39] V. Nair and G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *Proc. ICML*, 2010, pp. 1–8.
- [40] O. Russakovsky *et al.*, "ImageNet large scale visual recognition challenge," *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, Dec. 2015.
- [41] J. Donahue *et al.*, "DeCAF: A deep convolutional activation feature for generic visual recognition," in *Proc. 31st Int. Conf. Mach. Learn. (ICML)*, Beijing, China, Jun. 2014, pp. 21–26.
- [42] Y. Jia *et al.*, "Caffe: Convolutional architecture for fast feature embedding," in *Proc. 22nd ACM Int. Conf. Multimedia*, 2014, pp. 675–678.



Keze Wang received the B.S. degree in software engineering from Sun Yat-sen University, Guangzhou, China, in 2012. He is currently pursuing the Ph.D. degrees in computer science and technology with Sun Yat-sen University and Hong Kong Polytechnic University, Hong Kong, under the supervision of Prof. L. Lin and Prof. L. Zhang.

His current research interests include computer vision and machine learning.



Dongyu Zhang received the B.S. and Ph.D. degrees from the Harbin Institute of Technology, Harbin, China, in 2003 and 2010, respectively.

He is currently a Research Associate with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China. His current research interests include computer vision and machine learning.



Ya Li received the B.E. degree from Zhengzhou University, Zhengzhou, China, in 2002, the M.E. degree from Southwest Jiaotong University, Chengdu, China, in 2006, and the Ph.D. degree from Sun Yat-sen University, Guangzhou, China, in 2015.

She is currently a Lecturer with the School of Computer Science and Educational Software, Guangzhou University, Guangzhou. Her current research interests include computer vision and machine learning.



Ruimao Zhang received the B.E. degree from the School of Software, Sun Yat-sen University, Guangzhou, China, in 2011, where he is currently pursuing the Ph.D. degree in computer science with the School of Information Science and Technology.

He was a Visiting Ph.D. Student with the Department of Computing, Hong Kong Polytechnic University, Hong Kong, from 2013 to 2014. His current research interests include computer vision, pattern recognition, machine learning, and related applications.



Liang Lin (SM'14) received the B.S. and Ph.D. degrees from the Beijing Institute of Technology, Beijing, China, in 1999 and 2008, respectively.

He was a Post-Doctoral Research Fellow with the Department of Statistics, University of California at Los Angeles, Los Angeles, CA, USA, from 2008 to 2010. He was a Visiting Scholar with the Department of Computing, Hong Kong Polytechnic University, Hong Kong, and with the Department of Electronic Engineering, Chinese University of Hong Kong, Hong Kong. He is currently a Professor with the School of Computer Science, Sun Yat-sen University, Guangzhou, China. He has authored over 100 papers in top tier academic journals and conferences. His current research interests include new models, algorithms and systems for intelligent processing, and understanding of visual data, such as images and videos.

Prof. Lin received the Best Paper Runners-Up Award in ACM NPAR 2010, the Google Faculty Award in 2012, the Best Student Paper Award in the IEEE ICME 2014, and the Hong Kong Scholars Award in 2014. He currently serves as an Associate Editor of the IEEE TRANSACTIONS ON HUMAN-MACHINE SYSTEMS.