

# Embedding Visual Hierarchy with Deep Networks for Large-Scale Visual Recognition

Tianyi Zhao, Baopeng Zhang, Wei Zhang, Ning Zhou, Jun Yu, Jianping Fan

**Abstract**—In this paper, a level-wise mixture model (LMM) is developed by embedding visual hierarchy with deep networks to support large-scale visual recognition (i.e., recognizing thousands or even tens of thousands of object classes), and a Bayesian approach is used to adapt a pre-trained visual hierarchy automatically to the improvements of deep features (that are used for image and object class representation) when more representative deep networks are learned along the time. Our LMM model can provide an end-to-end approach for jointly learning: (a) the deep networks to extract more discriminative deep features for image and object class representation; (b) the tree classifier for recognizing large numbers of object classes hierarchically; and (c) the visual hierarchy adaptation for achieving more accurate indexing of large numbers of object classes hierarchically. By supporting joint learning of the tree classifier, the deep networks and the visual hierarchy adaptation, our LMM algorithm can provide an effective approach for controlling inter-level error propagation effectively, thus it can achieve better accuracy rates on large-scale visual recognition. Our experiments are carried on ImageNet1K and ImageNet10K image sets, and our LMM algorithm can achieve very competitive results on both the accuracy rates and the computation efficiency as compared with the baseline methods.

**Index Terms**—Large-scale visual recognition, level-wise mixture model (LMM), visual hierarchy adaptation, deep networks, tree classifier, Bayesian approach, object-group assignment matrix (group-object correlation matrix).

## I. INTRODUCTION

BY breaking the complex issue of feature learning into a set of small tasks hierarchically, deep learning [1], [2], [3], [4], [5] has demonstrated a divide-and-conquer process to learn more discriminative representations for large-scale visual recognition application: each neuron on the deep networks handles one small piece of the complex task for feature learning, and all these neurons can seamlessly collaborate to accomplish the complex task for feature learning in a coarse-to-fine fashion. For large-scale visual recognition application (i.e., recognizing thousands or even tens of thousands of object classes) [6], [7], [8], [9], [10], [11], [12], [13], the deep networks are usually trained in a supervised way by minimizing a flat loss function (such as cross-entropy). Some researchers have found that the neurons on the earlier layers of the deep networks are more 'common' but the neurons on the later layers are more 'specific' [14].

Tianyi Zhao, Baopeng Zhang, Wei Zhang, Ning Zhou, Jun Yu and Jianping Fan are with the Department of Computer Science, University of North Carolina, Charlotte, NC 28223, USA. e-mail: jfan@uncc.edu

This research is supported by National Science Foundation under 1651166-CNS.

Copyright (c) 2016 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending an email to pubs-permissions@ieee.org.

Even deep learning has achieved outstanding performances for many computer vision tasks, it still has room to improve: (1) strong inter-class visual similarities are typical in the domain of large-scale visual recognition especially when some object classes are fine-grained (visually-similar) [15], [16], [17], [18], but the  $N$ -way flat softmax classifier completely ignores the inter-task correlations; (2) ignoring the inter-task correlations completely may push the deep learning process away from the global optimum because the gradients of the joint objective function are not uniform for all the object classes, and such deep learning process may distract on discerning some particular object classes that are typically hard to be discriminated.

Another divide-and-conquer approach for supporting large-scale visual recognition is to leverage a pre-defined tree structure (visual hierarchy or concept ontology) [19], [20], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30], [31], [32], [33] to organize large numbers of object classes hierarchically. By training the tree classifier over a pre-defined tree structure hierarchically, the hierarchical visual recognition approach [34], [35], [36], [37] can provide multiple advantages: (a) Making coarse-to-fine predictions along a pre-defined tree structure can effectively rule out unlikely groups of object classes (i.e., irrelevant high-level nodes on the tree structure) at an early stage, thus it can achieve sub-linear computational complexity at test time; (b) For a given group (a high-level non-leaf node on the tree structure), the learning tasks for training the inter-related classifiers for its belonging fine-grained (visually-similar) object classes are strongly inter-related, thus multi-task learning can be used to train such inter-related classifiers jointly for enhancing their discrimination power [31], [32]; (c) For a given object class, the negative images for classifier training can be selected locally from other visually-similar object classes in the same group, and the issue of huge sample imbalance can be avoided effectively; (d) For a given group, its task space for object recognition is much smaller and uniform (i.e., distinguishing only a small number of fine-grained (visually-similar) object classes in the same group rather than separating all  $N$  object classes simultaneously [33]).

Based on these observations, it is very nature for us to ask ourselves the following question: *how can we integrate these two divide-and-conquer approaches (deep learning and hierarchical visual recognition) and benefit from both of them to exploit better solutions for large-scale visual recognition?*

In this paper, as shown in Fig. 1, a level-wise mixture model (LMM) is developed by using a tree classifier to replace traditional  $N$ -way flat softmax classifier in the deep networks,

where a visual hierarchy is embedded to: (a) organize large numbers of object classes hierarchically according to their inter-class visual similarities; (b) guide the process for joint learning of deep networks and tree classifier to make more effective splittings of the complex tasks for feature learning and hierarchical visual recognition. By leveraging group-object correlations (that are intuitively characterized by the visual hierarchy) to guide the process for jointly learning the deep networks and the tree classifier, our LMM model can boost the performance of large-scale visual recognition significantly and extract more robust and transferable features from the deep networks for image and object class representation. Because the deep features (outputs of the deep networks) are seamlessly integrated to learn the deep representations for large numbers of object classes and their inter-class visual similarities (that are used for constructing the visual hierarchy), the visual hierarchy should be adapted automatically when more representative deep networks are learned along the time, but it could be very expensive to reconstruct the visual hierarchy repeatedly. Based on this understanding, a Bayesian approach is further developed to effectively adapt the visual hierarchy during the end-to-end process for jointly learning the deep networks and the tree classifier.

In a summary, this paper has made the following *contributions*: (1) A level-wise mixture model (LMM) is developed to embed the visual hierarchy with the deep networks, so that we can leverage the group-object (inter-level) correlations (that are intuitively characterized by the visual hierarchy) to learn more representative deep networks and more discriminative tree classifier for supporting hierarchical visual recognition; (2) A Bayesian approach is developed to adapt the visual hierarchy to the improvements of deep class representations, e.g., learning more representative deep networks along the time may result in the improvements of the deep representations for large numbers of object classes and their inter-class visual similarities. Thus our LMM model can provide an end-to-end approach for jointly learning: (a) the deep networks to extract more discriminative features for image and object class representation; (b) the tree classifier (LMM model) for recognizing large numbers of object classes hierarchically; and (c) the visual hierarchy adaptation for achieving more accurate and hierarchical indexing of large numbers of object classes and identifying the tasks with similar learning complexities automatically. By supporting joint learning of the tree classifier, the deep networks and the visual hierarchy adaptation, our LMM algorithm can provide an effective approach for controlling inter-level error propagation effectively (i.e., inter-level error propagation is a critical issue for supporting hierarchical visual recognition[32], [33]). Our proposed algorithms have achieved very competitive results on ImageNet1K and ImageNet10K image sets as compared with the baseline methods.

The rest of this paper is organized as follows. Section II gives a brief review of the related works on deep learning, hierarchical visual recognition and tree structures. Section III introduces our level-wise mixture model (LMM) for learning the deep networks and the tree classifier jointly. Section IV describes our algorithm for visual hierarchy construction and adaptation. Section V provides our algorithm for learning the

deep networks, the tree classifier (LMM model) and the visual hierarchy adaptation jointly in an end-to-end fashion. Our experimental results for algorithm evaluation are presented in Section VI, and we conclude the paper and discuss the future works in Section VII.

## II. RELATED WORK

Deep learning[1], [2], [3], [4], [5] has demonstrated its outstanding abilities on learning more discriminative features and boosting the accuracy rates for large-scale visual recognition significantly. By learning more representative features and a  $N$ -way flat softmax classifier in an end-to-end fashion, most existing deep learning schemes have made one hidden assumption: the tasks for recognizing all the object classes are independent and share similar learning complexities. However, such assumption may not be true in many real-world applications, e.g., strong inter-class visual similarities are typical in the domain of large-scale visual recognition especially when some object classes are fine-grained (visually-similar) [15], [16], [17], [18], but the  $N$ -way flat softmax classifier completely ignores the inter-task correlations. Ignoring the inter-task correlations completely may push the deep learning process away from the global optimum because the gradients of the joint objective function are not uniform for all the object classes, especially when they have different inter-class visual similarities and learning complexities, as a result, the deep learning process may distract on discerning some particular object classes that are typically hard to be discriminated.

For large-scale visual recognition application (i.e., some object classes could have strong inter-class visual correlations), it is very attractive to develop new algorithms to deal with the issue of huge diversity of learning complexities more effectively, so that our deep learning schemes can effectively accomplish the task of learning more discriminative deep representations for distinguishing visually-similar object classes effectively. A few attempts have recently been made to exploit the tree structures (both concept ontology and visual hierarchy) in the deep learning models[38], [39], [5], [40], [33]. By integrating deep learning with multi-task learning, deep multi-task learning have received many attentions recently by using the deep networks to learn more representative features and integrating multi-task learning tools to learn inter-related classifiers jointly for separating such fine-grained (visually-similar) object classes more effectively [41], [42], [43], [44], [45], [46], [47], [48], [49], [50]. Most existing deep multi-task learning techniques assume that all the tasks are equally related and they may completely ignore the significant differences on the inter-task relationships (inter-class visual similarities) among large numbers of object classes [51], [52], [53], [53], [33], e.g., some object classes may have much stronger inter-class visual similarities and not all of them have same strengths on their inter-class visual similarities.

One intuitive way for exploiting the inter-task relationships (inter-class visual similarities) is to integrate a tree structure to organize large numbers of object classes hierarchically, e.g., the tasks for training the classifiers for the fine-grained (visually-similar) object classes under the same parent node

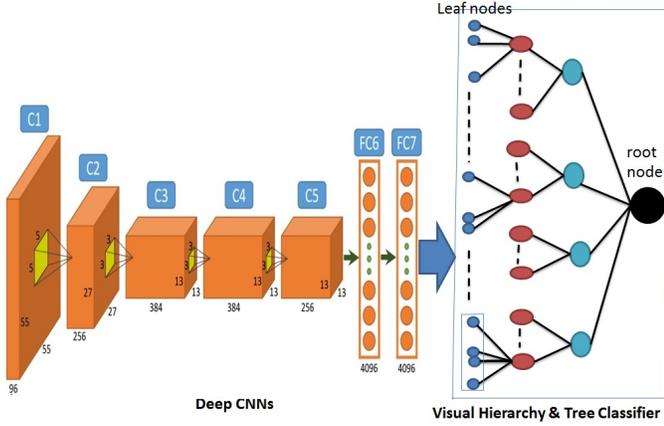


Fig. 1. The flowchart for embedding deep networks with visual hierarchy, where the tree classifier over the visual hierarchy is used to replace the traditional  $N$ -way flat softmax classifier.

(in the same group) may have stronger inter-task relationships and share similar learning complexities. Such tree structures can be categorized into two types: (a) concept ontology[19], [20], [21], [31], [30]; and (b) label tree or visual hierarchy[22], [23], [24], [25], [26], [27], [28], [29], [32], [33]. It is worth noting that the feature space is the common space for classifier training and visual recognition [54], e.g., both classifier training and visual recognition are performed in the feature space rather than in the semantic label space. Thus it could be more attractive to organize large numbers of object classes hierarchically in the feature space according to their inter-class visual correlations.

By integrating a tree structure to organize large numbers of object classes hierarchically and supervise the hierarchical process for tree classifier training, the hierarchical visual recognition approach [55], [34], [35], [36], [37], [30], [40], [31], [32] can provide many advantages, but it may seriously suffer from the problem of *inter-level error propagation*: the mistakes for the parent nodes will propagate to their child nodes until the leaf nodes [30], [32]. In addition, most existing approaches for hierarchical visual recognition focus on leveraging hand-crafted features for tree classifier training, thus it is very attractive to invest how deep features can be leveraged to improve hierarchical visual recognition [40], [33].

Most existing approaches for hierarchical visual recognition are static, but the process for joint learning of deep networks and tree classifier for large-scale visual recognition application is open-ended and dynamic: the deep networks for image and object class representation and the tree classifier for large-scale visual recognition may be improved along the time, e.g., more representative deep networks and more discriminative tree classifier may be achieved when more training images are added and back-propagation operations [56] are continuously performed to fine-tune the weights of the deep networks. Thus most existing approaches for hierarchical visual recognition may seriously suffer from the following problem: how to adapt the pre-trained tree structure (such as visual hierarchy) to the improvements of deep networks along the time? It is worth noting that the deep networks are used to obtain the

deep representations for large numbers of object classes and their inter-class visual similarities that are used for visual hierarchy construction. Thus it is very attractive to develop new approaches for jointly learning the deep networks, the tree classifier and the visual hierarchy adaptation in an end-to-end fashion.

### III. LEVEL-WISED MIXTURE MODEL (LMM)

Given  $N$  object classes being recognized, when a visual hierarchy is pre-trained for organizing  $N$  object classes hierarchically according to their inter-class visual similarities [32], [33], each level of the visual hierarchy can be treated as one particular partitioning  $\mathbb{T}_l$  of all these  $N$  object classes (i.e., assigning  $N$  object classes into a set of groups  $\mathbb{T}_l$  at the  $l$ th level of the visual hierarchy), followed by the distribution  $P: X \rightarrow \mathbb{T}_l$ ,  $X$  is the deep feature space for the training image set  $S$ ,  $X = h(S, u)$ ,  $u$  is the set of weights in the deep networks,  $h(\cdot)$  represents the mapping function of the deep networks,  $N_l$  is the number of groups at the  $l$ th level of the visual hierarchy, the distribution can be computed by the last layer of deep neural network, for example Softmax layer.

For a given group with the label  $t$  at the  $l$ th level of the visual hierarchy which contains  $C_t$  object classes, the probability  $P(y|t)$  for each of its  $C_t$  object classes can simply be defined as:  $P(y|t) = 1/C_t$ , e.g., we assume that the probability  $P(t|l, x, w_t^l)$  for the given group  $t$  (at the  $l$ th level of the visual hierarchy) is equally distributed among all its  $C_t$  object classes. Based on this assumption, for all the groups at the  $l$ th level of the visual hierarchy, the distribution over  $N$  object classes is defined as:

$$P(y|l, x, w^l) = \sum_t^{\mathbb{T}_l} I(t)P(t|l, x, w_t^l)P(y|t) \quad (1)$$

where  $w_l$  is the set of parameters for the node classifiers at the  $l$ th level of the visual hierarchy,  $w^l = \{w_t^l | t \in \mathbb{T}_l\}$ ,  $\mathbb{T}_l$  is the partitioning of  $N$  object classes at the  $l$ th level of the visual hierarchy, each layer classifier can be treated as an additional softmax layer over the deep networks,  $I(t)$  is an indication function and it is true when the object class with the label  $y$  belongs to the group with the label  $t$ ,  $P(t|l, x, w_t^l)$  is the distribution of the group  $t$  in one particular partitioning  $\mathbb{T}_l$  at the  $l$ th level of the visual hierarchy with deep representation  $x$ . It's worth noting that, we merging the deep neural network and the the Bayesian based Layer-wise Mixture Model by computing the probability  $P(t|l, x, w_t^l)$  by deep network.

By introducing a latent variable  $\theta$  to characterize the prior distribution over all the levels of the visual hierarchy, as shown in Fig. 2,  $l \sim \text{Cat}(\theta)$ ,  $l \in \{1, \dots, L\}$ , the mixture model  $P(y|x, W)$  is defined for modeling the probability of the object class with the label  $y$  given deep representation  $x$ :

$$P(y|x, W) = \sum_{l=0}^L \theta_l P(y|l, x, w^l) \quad (2)$$

where  $W$  is the set of parameters for all the node classifiers at different levels of the visual hierarchy,  $W = \{w^l | l \in \{1, \dots, L\}\}$ ,  $\theta_l$  is the leveraging parameter that is used to measure the contributions or effects from the node classifiers

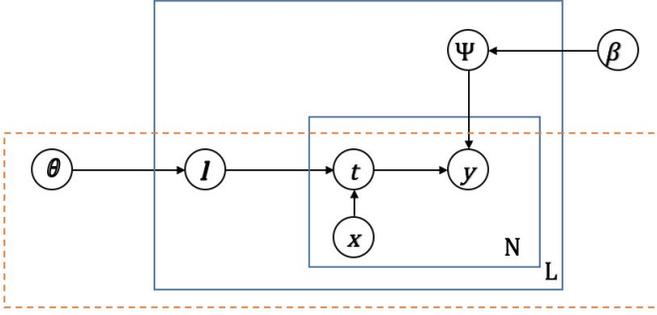


Fig. 2. Graph model for modeling the latent relationships between the object class  $y$  and the group  $t$ :  $l \sim \text{Cat}(\theta)$ ,  $t \sim \text{DL}(l, x)$ ,  $\Psi_t \sim \text{Dir}(\beta)$ ,  $y \sim \text{Cat}(\Psi_t)$ .

at the  $l$ th level of the visual hierarchy,  $L$  is the depth of the visual hierarchy, e.g., the total number of levels from the root node (which contains all these  $N$  object classes) to the leaf node (which contains only one particular object class).

By using our LMM model (tree classifier over the visual hierarchy) to replace traditional  $N$ -way flat softmax classifier, our deep networks for hierarchical visual recognition are shown in Fig. 1.

There are *two significant differences* between our deep networks and traditional deep CNNs [1]: (a) the tree classifier (LMM model) is used to replace the  $N$ -way flat softmax classifier, e.g., the tree classifier is defined as a set of node classifiers at different levels of the visual hierarchy; and (b) the group-class correlations (inter-level correlations) are leveraged to guide the process for jointly learning the deep networks and the tree classifier. Such group-object correlations (object-group assignments) are initially determined by a pre-trained visual hierarchy, and an object-group assignment matrix  $\Psi$  is further learned to measure such group-object (inter-level) correlations effectively (see Section IV). For a given group  $t$  at the  $l$ th level of the visual hierarchy, a softmax output is used to model the probability  $P(t|l, x, w_t^l)$  for the object class with the deep representation  $x$  to be assigned into the given group  $t$ .

After the deep networks and the tree classifier are learned jointly, for a given test image, it first goes through our deep networks to obtain its deep representation, and then such deep representation for the given test image goes through our tree classifier (our LMM model over the visual hierarchy) to obtain the prediction of its best-matching group  $t$  at each level of the visual hierarchy. The outputs from all these relevant node classifiers at different levels of the visual hierarchy are seamlessly integrated to produce the final prediction of the best-matching object class (at the bottom level of the visual hierarchy) for the given test image.

#### IV. VISUAL HIERARCHY ADAPTATION

To learn a pre-trained visual hierarchy for organizing  $N$  object classes hierarchically, we first need to calculate their inter-class visual similarities [32], [33]. For each object class, the deep networks are used to determine its deep representation from all its relevant images. The deep representation for each object class is simply obtained by averaging the deep features for all its relevant images [32], [33]. When such deep class



Fig. 3. One of our results on learning the pre-trained visual hierarchy for ImageNet1K image set.

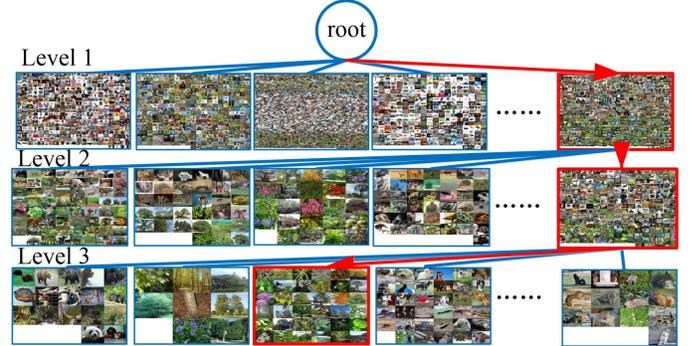


Fig. 4. One of our results on learning the pre-trained visual hierarchy for ImageNet10K image set.

representations are available for all these  $N$  object classes, we can further calculate a  $N \times N$  inter-class similarity matrix  $\mathbf{S}$  and its component  $S_{ij} = S(c_i, c_j)$  is defined as:

$$S_{i,j} = S(c_i, c_j) = \exp\left(-\frac{d(x_i, x_j)}{\sigma}\right) \quad (3)$$

where  $S(c_i, c_j)$  is the inter-class visual similarity between two object classes  $c_i$  and  $c_j$ ,  $d(x_i, x_j)$  is the distance between the deep class representations  $x_i$  and  $x_j$  for two object classes  $c_i$  and  $c_j$ ,  $\sigma$  is automatically determined by a self-tune technique. When the  $N \times N$  inter-class similarity matrix  $\mathbf{S}$  is available, hierarchical clustering is performed to learn the pre-trained visual hierarchy [32], [33]. Two of our experimental results on learning the pre-trained visual hierarchy are shown in Fig. 3 and Fig. 4.

In the previous section, all the object classes are treated equally, e.g., we assume that all the object classes have been assigned into the best-matching groups correctly and thus they have similar values (strengths) of the group-object correlations (i.e., all of them have good object-group assignments). In the real-world situation, it could be very hard for the pre-trained visual hierarchy to achieve optimal assignments for all the object classes because of two reasons: (a) The deep features for object class representation may not be discriminative enough because the deep networks for feature learning could be under-trained at the beginning, e.g., the deep networks may be improved along the time and the deep representations for large numbers of object classes and their inter-class visual similarities may also be improved along the time and thus the visual hierarchy should be adapted automatically (i.e., the pre-trained visual hierarchy should be adapted to the

improvements of deep networks or deep class representations); (b) The hierarchical clustering technique, which is used for visual hierarchy learning, may make mistakes.

As a result, some object classes may have stronger group-object correlations (i.e., they are strongly related with their belonging groups) because they are assigned into their best-matching groups correctly. On the other hand, some object classes may have weaker group-object correlations (i.e., they are weakly related with their belonging groups) because these object classes are assigned into some groups incorrectly. Thus it is very attractive to develop algorithms for adapting the pre-trained visual hierarchy automatically to the improvements of the outputs of the deep networks (i.e., deep features for object class representation), e.g., learning more representative deep networks along the time may result in both the improvements of deep class representations and the changes of their inter-class similarity matrix  $\mathbf{S}$  for visual hierarchy construction, which may further require the changes of the group-object correlations (e.g., the underlying object-group assignments should be improved adaptively).

In this work, a Bayesian approach is developed to estimate such group-object correlations accurately and adapt the pre-trained visual hierarchy to the improvements of the group-object correlations when more representative deep networks are learned along the time. We use an object-group assignment matrix  $\Psi$  to interpret the correspondences between the partitioning  $\mathbb{T}_l$  of all the object classes at the  $l$ th level of the visual hierarchy and the labels for all its belonging object classes,  $\Psi \in R^{N_l \times N}$ .  $\Psi_{t,y}$  represents the probability of the object class with the label  $y$  to be assigned into the given group with the label  $t$ , noted as  $P(y|t)$ . Given a pre-trained visual hierarchy, such object-group assignment matrix  $\Psi$  can automatically be determined and it can be initialed as described in Section III.

Our aim is to adapt the object-group assignment matrix  $\Psi$  automatically to the improvements of the outputs of the deep networks along the time. The method is to sample the group label  $t$  for each image, then update  $\Psi$  according to the sampling results. The probability model for characterizing such object-group assignments is described in Fig. 2.

In this paper, Dirichlet-categorical model is used,  $\Psi_t \sim Dir(\beta)$ ,  $y \sim Cat(\Psi_t)$ . The posterior likelihood of the group label  $t$  is obtained by integrating over the object-group assignment matrix  $\Psi$ :

$$P(t_i|T^{-i}, Y, X, l, \beta) \propto \frac{\Omega_{t_i, y_i}^{-i} + \beta_{y_i}}{\sum_y \Omega_{t_i, y}^{-i} + \beta_0} \times P(t_i|l, x_i, w_t^l) \quad (4)$$

where  $t_i$  is the predicted group label for the image  $i$ ,  $y_i$  is the label of the object class for the image  $i$ ,  $\mathbb{Y}$  is the set of all potential labels for the object classes,  $\{1, \dots, N\}$ ,  $Y$  is the set of the labels for all the images in the training set,  $T^{-i}$  is the set of group labels for all the images in the training set except the label for the  $i$ th image,  $\beta_0 = \sum_y \beta_y$ ,  $\Omega$  is a count matrix,  $\Omega_{t,y}$  is the number of images with the group label  $t$  and the object class label  $y$  as defined in Eq.(5),  $\Omega^{-i}$  is a count matrix without counting the  $i$ th image.

$$\Omega_{t,y} = \sum_{i=1}^{\Xi} I(t_i = t \wedge y_i = y) \quad (5)$$

where  $\Xi$  is the total number of labeled training images.

Eq.(4) has two components: (a) The first component is proportional to the number of images in the object class  $y_i$  which are assigned into the group  $t$ , e.g., more images in the group  $t$  come from the object class  $y_i$ , higher possibility for the images from the object class  $y_i$  to be assigned into the group  $t$ ; (b) The second component is the prediction from the deep networks, e.g., it gives the probability for the image  $i$  with the deep representation  $x_i$  to be assigned into the group  $t$ . The deep representation  $x_i$  for the image  $i$  may be improved during the iterative process for joint learning of deep networks and our LMM model (tree classifier). New group assignments for the object classes are obtained according to both the improvements of the deep class representations and the improvements of the inter-class similarity relationships for visual hierarchy construction.

The group-object correlation matrix (object-group assignment matrix)  $\Psi$  can be estimated automatically and its component  $\Psi_{t_i, y_i}$  is obtained as:

$$\Psi_{t_i, y_i} = \frac{\Omega_{t_i, y_i} + \beta_{y_i}}{\sum_y \Omega_{t_i, y} + \beta_0} \quad (6)$$

## V. JOINT LEARNING OF TREE CLASSIFIER, DEEP NETWORKS & VISUAL HIERARCHY ADAPTATION

In this work, an end-to-end approach is developed to jointly learn: (a) the LMM model (the tree classifier over the visual hierarchy) for recognizing large numbers of object classes hierarchically; (b) the group-object correlation matrix (object-group assignment matrix)  $\Psi$  for visual hierarchy adaptation; and (c) the deep networks to extract more discriminative features for image and object class representation. Our joint learning algorithm is illustrated in Algorithm 1.

### A. End-to-end Learning of LMM Model and Deep Networks

Joint learning of the LMM model and the deep networks is treated as an issue of Maximum Likelihood Estimation (MLE). The hierarchical recognition error for training the tree classifier over the visual hierarchy can be minimized by maximizing the log likelihood. The set of parameters for the node classifiers at different levels of the visual hierarchy is denoted as  $W$ , e.g., the tree classifier is represented as a set of node classifiers at different levels of the visual hierarchy. The weights for the deep networks are denoted as  $u$ , which are used to extract the feature vector  $x$  for image and object class representation. The log likelihood as defined in Eq.(7) is used for training the tree classifier over the visual hierarchy:

$$\mathcal{L}(W, u) = \log \left( \sum_{l=1}^L \theta_l \sum_t P(t|l, x, w_t^l) P(y|t) \right) \quad (7)$$

where  $w_t^l$  is the parameter of the node classifier for the group  $t$  at the  $l$ th level of the visual hierarchy,  $W$  is the set of parameters for all the node classifiers at different levels of the visual hierarchy,  $W = \{w_t^l |_{l=1}^L\}$ ,  $w_l$  is the set of parameters for all the group classifiers at the  $l$ th level of the visual hierarchy,  $w_l = \{w_t^l, t \in \mathbb{T}_l\}$ ,  $\theta_l$  is the leveraging parameter to measure the contributions or effects from the node classifiers at the  $l$ th

level of the visual hierarchy,  $\mathbb{T}_l$  is used to note the particular partitioning of all the object classes at the  $l$ th level of the visual hierarchy (i.e.,  $N_l$  groups at the  $l$ th level of the visual hierarchy),  $L$  is the depth of the visual hierarchy.

Given the log likelihood (joint objective function)  $\mathcal{L}(W, u)$ , all the parameters  $W$  for the node classifiers at different levels of the visual hierarchy and all the parameters (weights)  $u$  for the deep networks are learned by using stochastic gradient descent (SGD) [57], [58], [59], [60]. By maximizing the log likelihood (joint objective function)  $\mathcal{L}(W, u)$  as defined in Eq.(7), our LMM algorithm can achieve a path-based approach for learning the tree classifier over the visual hierarchy, e.g., learning the set of parameters  $W$  for the tree classifier. By back-propagating the gradients  $\frac{\partial \mathcal{L}(W, u)}{\partial u}$  of the objective function to fine-tune the weights  $u$  of the deep networks, our LMM algorithm can provide an end-to-end approach for learning the tree classifier and the deep networks jointly. Through the back-propagation process, our LMM algorithm can: (a) advise some 'common' neurons on the deep networks to learn more discriminative features for supporting more effective separation of the group nodes at the  $l$ th level of the visual hierarchy; and (b) guide some 'specific' neurons to learn more discriminative feature for achieving more accurate recognition of the fine-grained object classes at the bottom level of the visual hierarchy.

By fixing the group-object correlation matrix  $\Psi$  (i.e., fixing the object-group alignments or fixing the tree structure of the visual hierarchy), all the parameters  $W$  for the node classifiers at different levels of the visual hierarchy and all the parameters (weights)  $u$  for the deep networks are learned by maximizing the log likelihood  $\mathcal{L}(W, u)$ . The parameters of the tree classifier  $W$  are updated effectively by back-propagating the gradients  $\frac{\partial \mathcal{L}(W, u)}{\partial W}$  over the relevant node classifiers on the visual hierarchy. The parameters (weights) for the deep networks are fine-tuned effectively by back-propagating the gradients  $\frac{\partial \mathcal{L}(W, u)}{\partial u} = \frac{\partial \mathcal{L}(W, u)}{\partial x} \frac{\partial x}{\partial u}$  of the objective function (log likelihood) over the deep networks. The computation efficiency of our LMM algorithms is very competitive compared with the traditional  $N$ -way flat softmax method as shown in Eq. (8).

$$O(d \sum_i b^i) = O(d \frac{(1 - b^{\log_b N})}{(1 - b)} + Nd) = O(Nd) \quad (8)$$

where  $N$  is used to note the total number of object classes,  $d$  is used to note the dimension of the feature vector,  $b^i$  is used to note the number of branches for visual hierarchy construction.

### B. Back-Propagation Process

Given the log likelihood (joint objective function)  $\mathcal{L}(W, u)$ , the SGD method [57], [58], [59], [60] is used to learn the tree classifier and the deep networks jointly: (a) The gradients of the log likelihood  $\frac{\partial \mathcal{L}(W, u)}{\partial W}$  are calculated and such gradients are used to update the set of the parameters  $W = \{w^l\}_{l=1}^L$  of the tree classifier; (b) The gradients of the log likelihood

**Algorithm 1** Our algorithm for jointly learning deep networks, tree classifier and visual hierarchy adaptation

**Data:** Image set  $S$  with  $\Xi$  images, Label set  $Y$ , Pre-trained visual hierarchy.

Initialize group-object correlation matrix  $\Psi$  by using the pre-trained visual hierarchy.

**for** epoch = 1, ..., max of iterations **do**

    update the set of parameters  $W, u$  for LMM model (tree classifier) and deep networks by maximizing Eq. (7) based on the back-propagated the gradients  $\frac{\partial \mathcal{L}(W, u)}{\partial W}$ ,  $\frac{\partial \mathcal{L}(W, u)}{\partial u}$

**end for**

**repeat**

**for**  $i = 1, \dots, \Xi$  **do**

        sample the group label for the image  $i$  by using Eq.

(4)

        update the group-object correlation matrix  $\Psi$  by using Eq.(6)

        update the set of parameters  $W, u$  for LMM model (tree classifier) and deep networks by maximizing Eq. (7) based on the back-propagated the gradients  $\frac{\partial \mathcal{L}(W, u)}{\partial W}$ ,  $\frac{\partial \mathcal{L}(W, u)}{\partial u}$

**end for**

**until** converge

$\frac{\partial \mathcal{L}(W, u)}{\partial u}$  are calculated and such gradients are back-propagated to fine-tune the weights  $u$  of the deep networks.

$$\frac{\partial \mathcal{L}(W, u)}{\partial u} = \frac{\partial \mathcal{L}(W, u)}{\partial x} \frac{\partial x}{\partial u} \quad (9)$$

By maximizing the log likelihood  $\mathcal{L}(W, u)$ , our LMM algorithm can optimally minimize the hierarchical recognition error effectively, and learn a reliable group-object correlation matrix  $\Psi$  to control the inter-level error propagation effectively.

In our LMM model, for each node at the  $l$ th level of the visual hierarchy, a softmax output is used to estimate the probability  $P(t|l, x, w_t^l)$  for the image with deep representation  $x$  to be assigned into the group  $t$  at the  $l$ th level of the visual hierarchy. For a given image-class (feature-label) pair  $(x_i, y_i)$ , its group label  $t_i$  is defined as the ancestor for its object class  $y_i$  on the visual hierarchy:

$$z^l = \text{Softmax}(w^l x_i + b^l) \quad (10)$$

where  $w^l$  is the set of the parameters for all the node classifiers at the  $l$ th level of the visual hierarchy,  $b^l$  is the set of biases.

1) *Path-based Training of Tree Classifier:* In our path-based approach (our LMM model) for tree classifier training, the prediction probability  $P(y|x, W)$  of the object class  $y$  for the given image with deep representation  $x$  is obtained as:

$$z = \sum_l \theta^l z^l \Psi^l \quad (11)$$

where  $\theta^l$  is the leveraging parameter for characterizing the contributions or effects from the node classifiers at the  $l$ th level of the visual hierarchy,  $\Psi^l$  is used to note the group-object correlation matrix  $\Psi$  at the  $l$ th level of the visual hierarchy.

The hierarchical loss function for path-based training of the tree classifier is defined as the negative of the log likelihood:

$$\mathcal{L}(W, u) = -\log z_{y_i} \quad (12)$$

The gradients of the softmax output for the node classifiers at the  $l$ th level of the visual hierarchy are obtained as:

$$\frac{\partial \mathcal{L}(W, u)}{\partial z_t^l} = \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} \frac{\partial z_{y_i}}{\partial z_t^l} = \theta^l \Psi_{t, y_i}^l \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} \quad (13)$$

where  $\Psi_{t, y_i}^l$  is used to note the group-object correlation between the group  $t$  at the  $l$ th level of the visual hierarchy and the object class with the label  $y_i$ . Such gradient  $\frac{\partial \mathcal{L}(W, u)}{\partial z_t^l}$  can be used to measure the effects on the improvements of the parameters of the node classifier for the group  $t$  at the  $l$ th level of the visual hierarchy, and the weights of the deep network which are contributed by the softmax outputs at the  $l$ th level of the visual hierarchy. In Eq. (13),  $\theta^l$  decides how much effect from the given image  $(x_i, y_i)$  can be added to the node classifiers at the  $l$ th level of the visual hierarchy.  $\Psi_{t, y_i}^l$  decides how much effect from the group-object correlation can be added to the node classifier for the group  $t$  at the  $l$ th level of the visual hierarchy.

In the higher level of the visual hierarchy, if there is no adaptation on the object-group assignment matrix  $\Psi$ , Eq. (13) can be simplified as Eq. (14).

$$\frac{\partial \mathcal{L}(W, u)}{\partial z_{t_i}^l} = \theta^l \Psi_{t_i, y_i}^l \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} = \frac{\theta^l}{C_{t_i}^l} \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} \quad (14)$$

where  $C_{t_i}^l$  is used to note the number of object classes for the group  $t$  at the  $l$ th level of the visual hierarchy when the given image  $i$  is assigned into  $t$ .

For the gradients derived from the softmax outputs of the node classifiers at the  $l$ th level of the visual hierarchy, back-propagation [56] is used to leverage such gradients to update the parameters of the node classifier for the current node, and thus the corresponding modification on the classifier parameter  $\Delta w_t^l$  for the current node (i.e., group  $t$  at the  $l$ th level of the visual hierarchy) is defined as:

$$\Delta w_t^l = \epsilon \frac{\partial \mathcal{L}(W, u)}{\partial w_t^l} = \epsilon \frac{\theta^l}{C_{t_i}^l} \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} \frac{\partial z_{t_i}^l}{\partial w_t^l} \quad (15)$$

where  $\epsilon$  is used to note the learning rate. Such gradients for the group  $t$  at the  $l$ th level of the visual hierarchy are further used to update the classifier parameters for the lower-level nodes until the most relevant leaf nodes, which treat the group  $t$  as their ancestor on the visual hierarchy.

The gradients of the softmax output for the node classifiers at the bottom level of the visual hierarchy (i.e., for the object classes) are obtained as:

$$\frac{\partial \mathcal{L}(W, u)}{\partial z_{t_i}^{l_{bottom}}} = \theta^{l_{bottom}} \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} \quad (16)$$

where  $\theta^{l_{bottom}}$  is the leveraging parameter to measure the effects or contributions from the node classifiers at the bottom level of the visual hierarchy,  $t_i^{l_{bottom}}$  is used to note the group at the bottom level of the visual hierarchy. At the bottom level

of the visual hierarchy, each object class is treated as one single group, thus the group label for the bottom level  $t_i^{l_{bottom}}$  is same the object class label  $y_i$  for the given image  $i$ .

Our path-based approach can control the inter-level error propagation effectively: (a) For a given group  $t$  at the  $l$ th level of the visual hierarchy, the gradients of its node classifier as defined in Eq. (13) are used to update both the classifier parameters for itself and the classifier parameters for the lower-level nodes until the most relevant leaf nodes, which treat the group  $t$  as their ancestor on the visual hierarchy; (b) For a given object class  $y$  at the bottom level of the visual hierarchy,  $\Psi_{t, y}^{l_{bottom}} = 1$ , if  $t = y$ , else  $\Psi_{t, y}^{l_{bottom}} = 0$ , the gradients of its node classifier as defined in Eq. (13) are used to update only the classifier parameters for itself.

2) *Fine-tuning Network Weights*: The gradients derived from all the softmax outputs at different levels of the tree classifier (our LMM model) are then integrated to update the weights of the deep networks as given Eq. (17).

$$\frac{\partial \mathcal{L}(W, u)}{\partial x_i} = \sum_l \left\{ \frac{\theta^l}{C_{t_i}^l} \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} \frac{\partial z_{t_i}^l}{\partial x_i} \right\} \quad (17)$$

where  $\sum_l \theta^l = 1$ ,  $\theta^l$  is the leveraging parameter to measure the effects or contributions from the node classifiers at the  $l$ th level of the visual hierarchy. If  $\theta^{l_{bottom}} = 1$ , our LMM model is exactly same as the traditional  $N$ -way flat softmax model.

Without using non-overlapping constraint but with adaptation on the group-object correlation matrix  $\Psi$ , the derivative from our overall model is shown in Eqs. (18), (19).

$$\Delta w_t^l = \epsilon \sum_{t'}^{\tau_l} \left\{ \theta^l \Psi_{t', y_i}^l \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} \frac{\partial z_{t'}^l}{\partial w_t^l} \right\} \quad (18)$$

$$\frac{\partial \mathcal{L}(W, u)}{\partial x_i} = \sum_l \sum_t^{\tau_l} \left\{ \theta^l \Psi_{t, y_i}^l \frac{\partial \mathcal{L}(W, u)}{\partial z_{y_i}} \frac{\partial z_t^l}{\partial x_i} \right\} \quad (19)$$

By adapting the group-object correlation matrix and integrating back-propagation to leverage the gradients of the log likelihood (the objective function for tree classifier training) to fine-tune the weights of our deep networks, our LMM algorithm can provide an unique process for jointly learning: (a) the deep networks for image and object class representation; (b) the tree classifier (LMM model) for recognizing large numbers of object classes hierarchically; and (c) the visual hierarchy adaptation for achieving more accurate and hierarchical indexing of large numbers of object classes. By jointly learning the visual hierarchy adaptation, the tree classifier and the deep networks in an end-to-end fashion, our LMM algorithm can provide an effective solution for controlling inter-level error propagation effectively and achieve better accuracy rates on large-scale visual recognition.

### C. Visual Hierarchy Adaptation

First, the group-level labels for the images and the group-object correlation matrix (object-group assignment matrix)  $\Psi$  are initialized by the pre-trained visual hierarchy as described in Section III. The group-object correlation matrix  $\Psi$

TABLE I

The performance comparison between our LMM algorithm and the baseline method. Notations: P: group-object correlation matrix  $\Psi$  modification, TC: tree structure updating, w/: with, w/o: without.

Method	Prediction error (Top k)			
	1	2	5	10
baseline	45.13%	33.29%	21.67%	15.21%
LMM	45.05%	33.30%	21.64%	15.08%
LMM w/P, w/o TC	44.64%	32.93%	<b>21.24%</b>	14.87%
LMM w/P,TC	<b>44.29%</b>	<b>32.65%</b>	21.26%	<b>14.86%</b>

is adapted to the improvements of the outputs of the deep networks as described in Section IV.

Without any constraint, our Bayesian approach may result in an overlapping tree structure (i.e., some uncertain object classes can be assigned into multiple groups simultaneously rather than one single group). In our LMM model, it is really simple to cut the overlapped branches by using a non-overlapping constraint:  $\Psi_{t,y} = 0$ , where  $\Psi_{t,y} \neq \max_t \Psi_{t,y}$ .

#### D. Deep Learning with Regularization

The process for joint learning of deep networks and tree classifier can also be treated as the process of maximizing a posteriori estimation (MAP). The prior for the set of the parameters  $W$  (for all the node classifiers at different levels of the visual hierarchy) is chosen as Gaussian distribution with diagonal covariance. For example,

$$w_i \sim \mathcal{N}\left(0, \frac{1}{\alpha} I\right), w_i \in W \quad (20)$$

To learn the deep networks and the tree classifier jointly, we can maximize the log posterior likelihood:

$$\begin{aligned} \max_W \{\mathcal{L}'(W, u)\} &= \max_W \{\log P(y|x, W) + \log P(W)\} \\ &= \max_W \left\{ \log P(y|x, W) - \frac{\alpha}{2} \|W\|^2 \right\} \end{aligned} \quad (21)$$

The first term in Eq.(21) is same as that in Eq.(7). The second term is L2-norm regularization over the classifier parameter  $W$ , which is used to control over-fitting and learn more discriminative tree classifier. Eq.(21) can be used to replace Eq.(7) in Algorithm1.

## VI. EXPERIMENTAL RESULTS FOR ALGORITHM EVALUATION

Our experiments are carried over two image sets: ImageNet1K and ImageNet10K. ImageNet1K image set contains 1,000 object classes, which are mutual exclusive or overlap but not subsumption in the semantic space. As shown in Fig. 3, the visual hierarchy is pre-trained to organize 1,000 object classes hierarchically. ImageNet10K image set contains 10,184 image categories, which come from different levels of the concept ontology and some of them could be subsumption of others, e.g., not all of these 10,184 image categories are semantically atomic (mutually exclusive) because some of them are from the high-level non-leaf nodes of the concept ontology and they are not semantically atomic with others. Thus the concept ontology is incorporated to decompose such

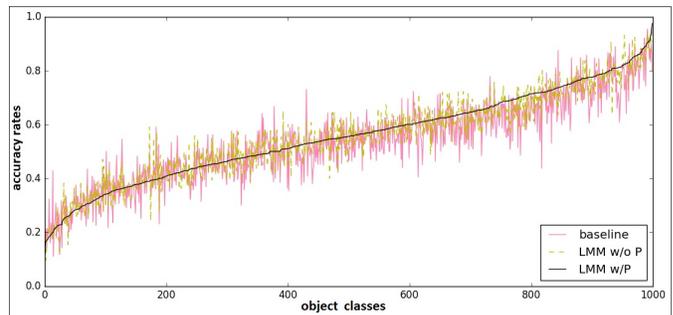


Fig. 5. The comparison on the accuracy rates for 1,000 object classes in ImageNet1K image set.

high-level image categories (from the non-leaf nodes of the concept ontology) into multiple object classes (at the leaf nodes of the concept ontology), and 7,756 object classes are finally identified for ImageNet10K image set. As shown in Fig. 4, the visual hierarchy is also pre-trained to organize 7,756 object classes hierarchically.

We have compared our LMM algorithm with the baseline method [1], [61] and our comparison experiments focus on evaluating multiple factors: (a) whether our LMM algorithm can control the inter-level error propagation more effectively as compared with other baseline methods and achieve better accuracy rates on large-scale visual recognition; (b) whether our LMM algorithm can jointly learn the tree classifier, the deep networks and the visual hierarchy adaptation effectively; and (c) whether our LMM algorithm can achieve higher prediction confidences.

Our experiments are implemented on Theano[62], with one GPU NVIDIA GTX 980i. In our experiments, the learning rate is set as 0.0001, the dropout rate is set as 0.5 to prevent over-fitting. We take an initial from Theano Alexnet pre-trained [61].

#### A. Experimental Results on Imagenet1K

For the ImageNet1K image set, a deep CNNs with a  $N$ -way flat softmax classifier is first trained as the baseline [1], [61]. Second, our LMM model (tree classifier) is trained by using a pre-trained visual hierarchy (without modifying the group-object correlation matrix  $\Psi$  along the time). Third, our LMM model is trained with modifications on the group-object correlation matrix  $\Psi$ . Finally, our LMM model is trained by jointly learning the deep networks, the tree classifier, and the visual hierarchy adaptation in an end-to-end fashion. The experimental results on the average accuracy rates are shown in Table I, one can easily observe that: (a) our LMM algorithm can successfully outperform the baseline method; (b) When our LMM algorithm jointly learns the deep networks, the tree classifier, and the visual hierarchy adaptation in an end-to-end fashion, it can achieve the best performance.

The comparisons on the accuracy rates are shown in Fig. 5, for all these 1,000 object classes in the ImageNet1K image set, 56% object classes have obtained better accuracy rates by using our LMM model, 19% object classes remain no obvious changes on their accuracy rates, 25% object classes loss more than 1% accuracy rates. Some visual recognition examples and

their confidence scores are shown in Fig. 6, Fig. 7 and Fig. 8, one can easily observe that our LMM algorithm can achieve more accurate recognition with higher confidence scores.

1) *Accuracy Rates vs. Learning Complexities*: As we mentioned before, one goal for embedding the visual hierarchy with the deep networks is to deal with the issues of strong inter-class visual similarities and diverse learning complexities for large-scale visual recognition. Thus it is very attractive to develop new algorithms to: (a) separate such visually-similar object classes with similar learning complexities from others; and (b) train the inter-related classifiers for such visually-similar object classes jointly. By integrating the visual hierarchy to assign the visually-similar object classes into the same group, our LMM algorithm can learn their inter-related classifiers jointly. Because such visually-similar object classes share similar learning complexities, the gradients of their joint objective function are more uniform and homogeneous, so that the back-propagation process can easily stick on reaching a global optimum effectively. As a result, our LMM algorithm can achieve higher accuracy rates on distinguishing such visually-similar object classes which are typically hard to be separated.

Because the deep networks and the tree classifier may not be discriminative enough at the beginning or the pre-trained visual hierarchy may make wrong assignments for some object classes, our LMM model may have low accuracy rates for these object classes which are initially assigned into wrong groups because their deep class representations are not discriminative enough at the beginning. After few iterations, the accuracy rates for these object classes can be improved significantly when they are finally re-assigned into their best-matching groups correctly. The reason for such improvements is that our LMM model can jointly learn the deep networks, the tree classifier and the visual hierarchy (i.e., the object-group assignment matrix), thus they can be improved simultaneously, e.g., more accurate deep class representations (more representative deep networks) can result in more accurate assignments of object classes (more accurate visual hierarchy) and learn more discriminative tree classifier. The names for some of those object classes with improved accuracy rates are listed in Table III.

In Fig 9, we illustrate our experimental results for 4 groups with different numbers of visually-similar object classes. From these experimental results, one can observe that there have good correspondences between the accuracy rates for object recognition and the strengths of the group-object correlations. From these experimental results, one can observe that: (a) The object classes, which have large values (strengths) of group-object correlations at the beginning, may have small improvements on their accuracy rates along the time because they have already been assigned into their best-matching groups correctly; (b) The object classes, which have low values (strengths) of group-object correlations at the beginning, may have big improvements on their accuracy rates along the time when more representative deep networks are learned and they are finally re-assigned into their best-matching groups correctly.

The reasons are that: (a) The object classes, which have

large values (strengths) of group-object correlations at the beginning, have already been assigned into their best-matching groups correctly by the pre-trained visual hierarchy, and the deep representations for those object classes have already been exploited for training the group-level classifiers effectively and thus they may have less contributions on improving the accuracy rates at the group level; (b) The object classes, which have low values (strengths) of group-object correlations at the beginning, may initially be assigned into wrong groups by the pre-trained visual hierarchy and can be re-assigned into their best-matching groups correctly because of visual hierarchy adaptation, and thus those re-assigned object classes may have more contributions on improving the group-level classification performance along the time. Thus the object classes, which have low values (strengths) of group-object correlations at the beginning, could have significant improvement on the accuracy rates at the object class level.

Overall, by leveraging the visual hierarchy to assign the visually-similar object classes with similar learning complexities into the same group and learn their inter-related classifiers jointly, the gradients of their joint objective function are more uniform and homogeneous, thus our LMM algorithm can obtain global optimum effectively and result in more discriminative tree classifier for large-scale visual recognition.

2) *Prediction Confidences*: For large-scale visual recognition application, it is also very important to evaluate the confidences of the predictions for object recognition. As shown in Fig. 6, Fig. 7 and Fig. 8, by assigning the visually-similar object classes with similar learning complexities into the same group and learning their inter-related classifiers jointly, our LMM model can obtain higher prediction confidence scores as compared with the baseline method.

3) *Visual Hierarchy Adaptation and Object-Group Reassignment*: The effects of visual hierarchy adaptation are evaluated by comparing multiple approaches: (1) our LMM model with an initial group-object correlation matrix  $\Psi$  (provided by a pre-trained visual hierarchy); (2) our LMM model with modification of group-object correlation matrix via visual hierarchy adaptation; and (3) our LMM model with modification of group-object correlation matrix and non-overlapping constraint. The names of object classes in that group are shown in Table II, the names of the object classes which are initially assigned into that group are listed in red background in Table II. The re-assignments of the object classes from our visual hierarchy adaptation method are listed in Table II. The names listed in the blue box are the new object classes which are re-assigned into this particular group and only 20% object classes are re-assigned because of visual hierarchy adaptation, thus the pre-trained visual hierarchy can achieve reasonable performance on assigning the visually-similar object classes with similar learning complexities into the same group [32], [33].

4) *Inter-Class Separability Enhancement*: By assigning the visually-similar atomic object classes with similar learning complexities into the same group, as shown in Fig. 10, our LMM algorithm can significantly enhance their inter-class separability by focusing on learning more discriminative deep representations and node classifiers to enlarge their inter-



TABLE III

The lists of object classes with increased and decreased accuracy rates when more representative deep networks are learned along the time.

	Accuracy Increased (>0.05)	Accuracy decreased (<0.0)
High correlation(>0.04)	great white shark, tiger shark, electric ray, stingray, brambling, goldfinch, junco, indigo bunting, robin, bulbul, green lizard, green snake, water snake, green mamba, sidewinder, ruffed grouse, partridge, lorikeet, hornbill, red-breasted merganser, wombat, black stork, little blue heron, European gallinule, red-backed sandpiper, dowitcher, oystercatcher, Rhodesian ridgeback, beagle, English foxhound, Ibizan hound, American Staffordshire terrier, Chesapeake Bay retriever, Appenzeller, bull mastiff, Great Dane, Saint Bernard, Eskimo dog, malamute, red wolf, dhole, grey fox, leopard, ground beetle, long-horned beetle, grasshopper, cockroach, monarch, marmot, mink, black-footed ferret, guenon, patas, langur, howler monkey, African elephant, bathtub, jinrikisha, mountain bike, police van, sarong, screen, snowmobile, submarine, toyshop, plate, lakeside	chickadee, bald eagle, American chameleon, vine snake, night snake, horned viper, wolf spider, tusker, spoonbill, crane, pelican, Chihuahua, basset, bloodhound, Walker hound, redbone, Saluki, English setter, Cardigan, tiger cat, lynx, leaf beetle, weevil, cricket, leafhopper, cabbage butterfly, beaver, weasel, polecat, macaque, colobus, altar, analog clock, bakery, bookshop, cab, cloak, confectionery, desktop computer, dining table, dock, dome, drilling platform, gown, grocery store, handkerchief, lawn mower, library, minivan, mitten, moped, palace, pickup, pillow, purse, quilt, rocking chair, screwdriver, shield, soup bowl, triumphal arch, tub, yawl, ice cream, potpie
Low correlation(<0.01)	European fire salamander, spotted salamander, leatherback turtle, mud turtle, peacock, fiddler crab, spiny lobster, hermit crab, king penguin, Bedlington terrier, Scotch terrier, Gordon setter, Irish water spaniel, affenpinscher, Pomeranian, keeshond, Brabancon griffon, sloth bear, hippopotamus, barber chair, baseball, chain saw, chime, forklift, French horn, guillotine, hourglass, maraca, ocarina, parking meter, pedestal, pinwheel, reflex camera, rotisserie, school bus, tennis ball, warplane, wing, volcano	great grey owl, bullfrog, Gila monster, triceratops, scorpion, centipede, rock crab, isopod, Irish setter, komondor, giant panda, assault rifle, backpack, bikini, bullet train, dumbbell, electric guitar, grille, honeycomb, mountain tent, oil filter, Polaroid camera, projectile, remote control, running shoe, sandal, sliding door, space bar, stethoscope, vacuum, cardoon, geyser, gyromitra

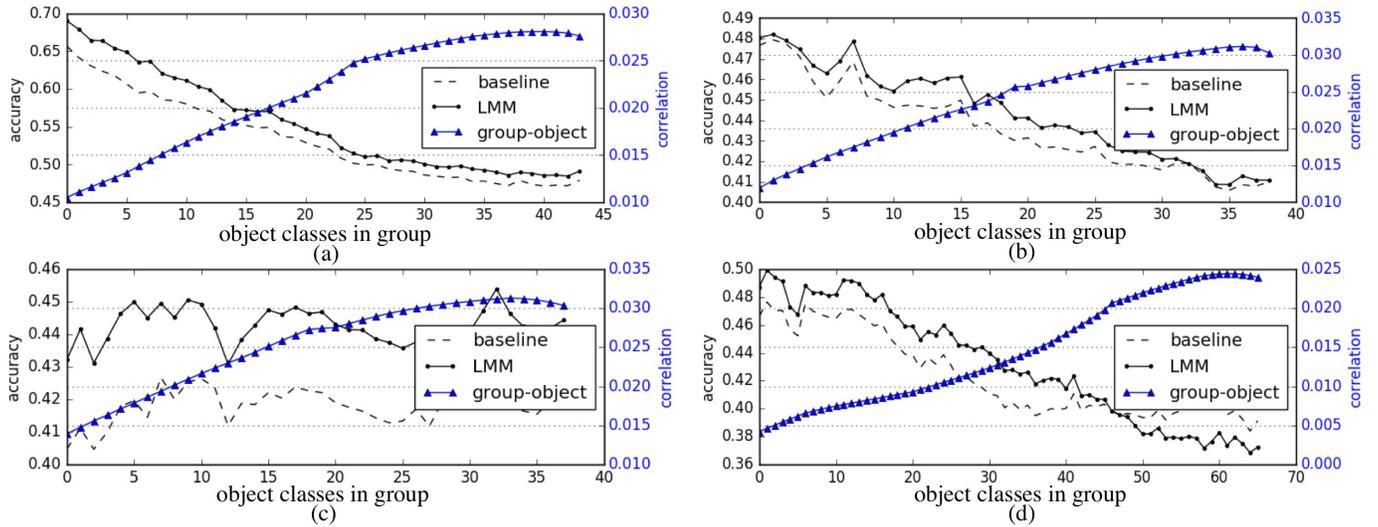


Fig. 9. The experimental results for 4 groups of object classes on their correspondences between the accuracy rates for object recognition and the strengths of the group-object correlations.

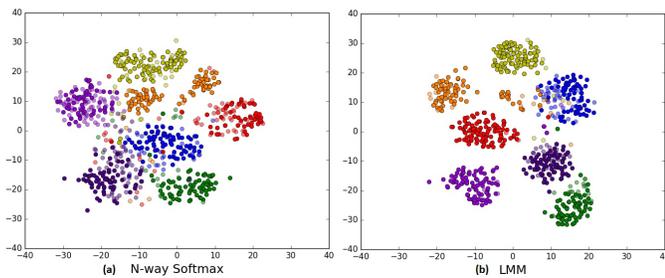


Fig. 10. The comparison on the inter-class separability for the visually-similar object classes in the same group.

class margins. By assigning the visually-similar atomic object classes with similar learning complexities into the same group, the gradients of their joint objective function are more uniform and homogeneous, thus our LMM algorithm can obtain global

optimum effectively and enlarge their inter-class margins significantly.

### B. Experimental Results on Imagenet10K

For the ImageNet10K image set, we take the parameters from the deep networks, which are already trained for the Imagenet1K image set, as the initials. We then use the images from ImageNet10K to train our LMM model. We have evaluated several kinds of initials, such as the baseline method [1], [61] and our LMM model. We finally find that taking the parameters in our LMM model (which is already trained for Imagenet1K image set) to initialize our deep networks for ImageNet10K image set can allow us to achieve the best performance.

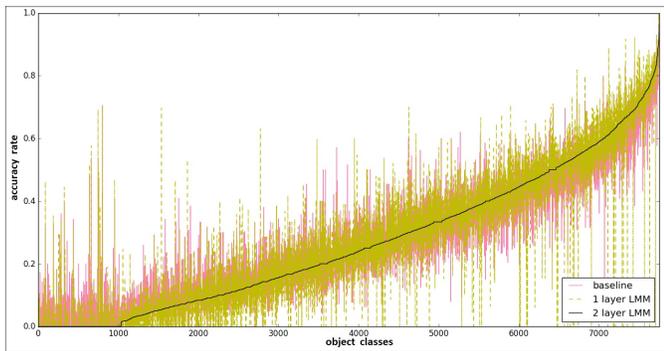


Fig. 11. The comparison on the accuracy rates for all the object classes in ImageNet10K image set.

TABLE IV

The performance comparison between multiple approaches over the ImageNet10K image set.

Method	Prediction error (Top k)			
	1	2	5	10
$N$ -way Softmax	70.30%	60.33%	47.99%	39.20%
one-level LMM	69.60%	<b>59.13%</b>	<b>46.87%</b>	38.82%
two-level LMM	<b>69.50%</b>	59.39%	46.88%	<b>38.07%</b>

Because the training cost for the ImageNet10K image set is very high, we only perform our LMM method over one visual hierarchy whose number of branch at each node is 100 (i.e., each parent node contains 100 sibling child nodes). We have also compared our LMM algorithm over the ImageNet10K image set under the following settings: (a)  $N$ -way flat softmax classifier [1], [61]. (b) One-level LMM model: each group contains only one object class. The group-object correlation matrix is same as the confusion matrix between all the object classes,  $\Psi \in R^{N \times N}$ . (c) Two-level LMM model: one level for coarse-grained groups and one level for the visually-similar object classes. The experimental results on average accuracy rates are listed on Table IV, one can observe that our LMM algorithm can successfully outperform the  $N$ -way flat softmax classifier. The comparisons on the accuracy rates are shown in Fig. 11, for all these object classes in the ImageNet10K image set, 68% object classes have obtained better accuracy rates by using our LMM model, 19% object classes remain no obvious changes on their accuracy rates, 13% object classes loss more than 1% accuracy rates.

By integrating the visual hierarchy to assign the visually-similar object classes into the same group, such visually-similar object classes in the same group may share similar learning complexities, thus the gradients of their joint objective function are more uniform and homogeneous, so that our LMM algorithm can easily stick on reaching a global optimum effectively and achieve higher accuracy rates on large-scale visual recognition.

## VII. CONCLUSIONS AND FUTURE WORKS

A level-wise mixture model (LMM) is developed in this paper to boost the performance of large-scale visual recognition. Our LMM model can provide an end-to-end approach to jointly learn the deep networks for image and object class

representation, the tree classifier for recognizing large numbers of object classes hierarchically and the visual hierarchy adaptation for achieving more accurate and hierarchical indexing of large numbers of object classes, thus our LMM algorithm can also provide an effective approach for controlling inter-level error propagation effectively and achieve better accuracy rates on large-scale visual recognition. By seamlessly integrating two divide-and-conquer approaches (deep learning and hierarchical visual recognition), we have found that these two approaches can benefit from each other to exploit better solutions for large-scale visual recognition. Our experimental results on ImageNet1K and ImageNet10K image sets have demonstrated that our LMM algorithm can achieve competitive results on the accuracy rates as compared with the baseline.

Many other deep networks, such as VGG [2], GoogleNets [3], Resnet [4], have successfully designed to recognize 1,000 object classes. Thus it is very attractive to leverage these successful designs of deep networks [2], [3], [4] to configure our deep networks and evaluate the performances of our LMM algorithm when different types of deep networks are used. Because these complex deep networks [2], [3], [4] have achieved better performance than AlexNet [1], [63], [64] used in this paper, we can expect that using these complex deep networks can allow our LMM algorithm to achieve higher accuracy rates on large-scale visual recognition.

## REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems* 25. Curran Associates, Inc., 2012, pp. 1097–1105.
- [2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," *CoRR*, vol. abs/1409.4842, 2014.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015.
- [5] T. Kong, A. Yao, Y. Chen, and F. Sun, "Hypernet: Towards accurate region proposal generation and joint object detection," *CoRR*, vol. abs/1604.00600, 2016.
- [6] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid, "Good Practice in Large-Scale Learning for Image Classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 3, pp. 507–520, Mar. 2014.
- [7] T. Dean, M. Ruzon, M. Segal, J. Shlens, S. Vijayanarasimhan, and J. Yagnik, "Fast, accurate detection of 100,000 object classes on a single machine," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, Washington, DC, USA, 2013.
- [8] J. Deng, A. C. Berg, K. Li, and L. Fei-Fei, "What does classifying more than 10,000 image categories tell us?" in *Proceedings of the 11th European Conference on Computer Vision: Part V*, ser. ECCV'10. Berlin, Heidelberg: Springer-Verlag, 2010, pp. 71–84.
- [9] E. P. Xing and B. Zhao, "Sparse output coding for large-scale visual recognition," vol. 00, pp. 3350–3357, 2013.
- [10] J. Deng, J. Krause, A. Berg, and L. Fei-Fei, "Hedging your bets: Optimizing accuracy-specificity trade-offs in large scale visual recognition," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Providence, RI, USA, June 2012.
- [11] Y. Lin, L. Cao, F. Lv, S. Zhu, M. Yang, T. Cour, K. Yu, and T. S. Huang, "Large-scale image classification: fast feature extraction and svm training," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2011.
- [12] J. Deng, N. Ding, Y. Jia, A. Frome, K. Murphy, S. Bengio, Y. Li, H. Neven, and H. Adam, "Large-scale object classification using label relation graphs," in *European Conference on Computer Vision*, ser. Lecture Notes in Computer Science, vol. 8689. Springer, 2014, pp. 48–64.

- [13] J. Li, S. Albaradei, Y. Wang, and L. Cao, "Learning mid-level features from object hierarchy for image classification," in *Proceedings of IEEE Winter Conference on Applications of Computer Vision*, 2014.
- [14] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" *CoRR*, vol. abs/1411.1792, 2014.
- [15] A. R. Sfar, N. Boujemaa, and D. Geman, "Vantage feature frames for fine-grained categorization," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, 2013, pp. 835–842.
- [16] R. Farrell, O. Oza, N. Zhang, V. I. Morariu, T. Darrell, and L. S. Davis, "Birdlets: Subordinate categorization using volumetric primitives and pose-normalized appearance," in *IEEE International Conference on Computer Vision, ICCV*, nov. 2011.
- [17] G. Martinez-Munoz, T. Dietterich, S. Todorovic, A. Yamamuro, W. Zhang, R. Paasch, A. Moldenke, N. Larios, D. Lytle, E. Mortensen, N. Payet, and L. Shapiro, "Dictionary-free categorization of very similar objects via stacked evidence trees," *2013 IEEE Conference on Computer Vision and Pattern Recognition*, vol. 00, no. undefined, pp. 549–556, 2009.
- [18] J. Deng, J. Krause, and L. Fei-Fei, "Fine-grained crowdsourcing for fine-grained recognition," *2013 IEEE Conference on Computer Vision and Pattern Recognition*, vol. 00, no. undefined, pp. 580–587, 2013.
- [19] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *CVPR09*, 2009.
- [20] A. Hauptmann, J. Curtis, L. Kennedy, J. R. Smith, S.-F. Chang, J. Tesic, W. Hsu, and M. Naphade, "Large-scale concept ontology for multimedia," *IEEE MultiMedia*, vol. 13, no. undefined, pp. 86–91, 2006.
- [21] M. Marszalek and C. Schmid, "Semantic hierarchies for visual object recognition," *Computer Vision and Pattern Recognition, 2007. CVPR '07. IEEE Conference on*, pp. 1–7, June 2007.
- [22] —, "Constructing category hierarchies for visual recognition," in *European Conference on Computer Vision*, ser. Lecture Notes in Computer Science, vol. 5305. Springer, 2008, pp. 479–491.
- [23] D. Nister and H. Stewenius, "Scalable recognition with a vocabulary tree," in *Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Volume 2*, ser. CVPR '06. Washington, DC, USA: IEEE Computer Society, 2006, pp. 2161–2168.
- [24] G. Griffin and P. Perona, "Learning and using taxonomies for fast visual categorization," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, 2008.
- [25] S. Bengio, J. Weston, and D. Grangier, "Label embedding trees for large multi-class tasks," in *Advances in Neural Information Processing Systems 23*. Curran Associates, Inc., 2010, pp. 163–171.
- [26] E. Bart, I. Porteous, P. Perona, and M. Welling, "Unsupervised learning of visual taxonomies," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008.
- [27] J. Deng, S. Satheesh, A. C. Berg, and F. Li, "Fast and balanced: Efficient label tree learning for large scale object recognition," in *Advances in Neural Information Processing Systems 24*. Curran Associates, Inc., 2011, pp. 567–575.
- [28] B. Liu, F. Sadeghi, M. F. Tappen, O. Shamir, and C. Liu, "Probabilistic label trees for efficient large scale image classification," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, 2013, pp. 843–850.
- [29] J. Sivic, B. C. Russell, A. Zisserman, W. T. Freeman, and A. A. Efros, "Unsupervised discovery of visual object class hierarchies," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [30] J. J. McAuley, A. Ramisa, and T. S. Caetano, "Optimization of robust loss functions for weakly-labeled image taxonomies: An imagenet case study," *International Journal of Computer Vision*, vol. 104, no. 3, pp. 343–361, 2013.
- [31] J. Fan, Y. Gao, and H. Luo, "Integrating concept ontology and multitask learning to achieve more effective classifier training for multilevel image annotation," *IEEE Trans. Image Processing*, vol. 17, no. 3, pp. 407–426, 2008.
- [32] J. Fan, N. Zhou, J. Peng, and L. Gao, "Hierarchical learning of tree classifiers for large-scale plant species identification," *IEEE Trans. Image Processing*, vol. 24, no. 11, pp. 4172–4184, 2015.
- [33] J. Fan, T. Zhao, Z. Kuang, Y. Zhang, J. Zhang, J. Yu, and J. Peng, "Hd-ntl: Hierarchical deep multi-task learning for large-scale visual recognition," *IEEE Trans. Image Processing*, vol. 26, 2017.
- [34] O. Dekel, J. Keshet, and Y. Singer, "Large margin hierarchical classification," in *Proceedings of the Twenty-first International Conference on Machine Learning*, ser. ICML '04. New York, NY, USA: ACM, 2004, pp. 27–.
- [35] M. W. Dengyong Zhou, Lin Xiao, "Hierarchical classification via orthogonal transfer," Tech. Rep., May 2011.
- [36] M. Sun, W. Huang, and S. Savarese, "Find the best path: an efficient and accurate classifier for image hierarchies," in *Proceedings of the International Conference on Computer Vision*, 2013.
- [37] J. Wang, X. Shen, and W. Pan, "On large margin hierarchical classification with multiple paths," *Journal of the American Statistical Association*, vol. 104, no. 487, pp. 1213–1223, 2009.
- [38] Z. Yan, V. Jagadeesh, D. DeCoste, W. Di, and R. Piramuthu, "HD-CNN: hierarchical deep convolutional neural network for image classification," *CoRR*, vol. abs/1410.0736, 2014.
- [39] N. Srivastava and R. R. Salakhutdinov, "Discriminative transfer learning with tree-based priors," in *Advances in Neural Information Processing Systems 26*. Curran Associates, Inc., 2013, pp. 2094–2102.
- [40] P. Kotschieder, M. Fiterau, A. Criminisi, and S. R. Bulò, "Deep neural decision forests," in *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, 2015, pp. 1467–1475.
- [41] Z. Zhang, P. Luo, C. C. Loy, and X. Tang, "Facial landmark detection by deep multi-task learning," in *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VI*, 2014, pp. 94–108.
- [42] Y. Huang, W. Wang, L. Wang, and T. Tan, "Multi-task deep neural network for multi-label learning," *2013 IEEE International Conference on Image Processing*, pp. 2897–2900, 2013.
- [43] Z. Wu, Y.-G. Jiang, J. Wang, J. Pu, and X. Xue, "Exploring inter-feature and inter-class relationships with deep neural networks for video classification," in *Proceedings of the 22Nd ACM International Conference on Multimedia*, ser. MM '14. New York, NY, USA: ACM, 2014, pp. 167–176.
- [44] P. Teterwak and L. Torresani, "Shared Roots: Regularizing Deep Neural Networks through Multitask Learning," Dartmouth College, Computer Science, Hanover, NH, Tech. Rep. TR2014-762, June 2014.
- [45] S. Li, Z. Liu, and A. B. Chan, "Heterogeneous multi-task learning for human pose estimation with deep convolutional neural network," *CoRR*, vol. abs/1406.3474, 2014.
- [46] X. Li, L. Zhao, L. Wei, M. Yang, F. Wu, Y. Zhuang, H. Ling, and J. Wang, "Deepsaliency: Multi-task deep neural network model for salient object detection," *CoRR*, vol. abs/1510.05484, 2015.
- [47] C. Zhang and Z. Zhang, "Improving multiview face detection with multi-task deep convolutional neural networks," *WACV*, 2014.
- [48] B. Hariharan, P. A. Arbeláez, R. B. Girshick, and J. Malik, "Hypercolumns for object segmentation and fine-grained localization," *CoRR*, vol. abs/1411.5752, 2014.
- [49] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," *CoRR*, vol. abs/1411.4038, 2014.
- [50] T. Kong, A. Yao, Y. Chen, and F. Sun, "Hypernet: Towards accurate region proposal generation and joint object detection," *CoRR*, vol. abs/1604.00600, 2016.
- [51] A. R. Sfar, N. Boujemaa, and D. Geman, "Vantage feature frames for fine-grained categorization," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, 2013, pp. 835–842.
- [52] R. Farrell, O. Oza, N. Zhang, V. I. Morariu, T. Darrell, and L. S. Davis, "Birdlets: Subordinate categorization using volumetric primitives and pose-normalized appearance," in *IEEE International Conference on Computer Vision, ICCV 2011, Barcelona, Spain, November 6-13, 2011*, 2011, pp. 161–168.
- [53] G. Martinez-Munoz, N. L. Delgado, E. N. Mortensen, W. Zhang, A. Yamamuro, R. Paasch, N. Payet, D. A. Lytle, L. G. Shapiro, S. Todorovic, A. Moldenke, and T. G. Dietterich, "Dictionary-free categorization of very similar objects via stacked evidence trees," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, 2009, pp. 549–556.
- [54] J. Fan, X. He, N. Zhou, J. Peng, and R. Jain, "Quantitative characterization of semantic gaps for learning complexity estimation and inference model selection," *IEEE Trans. on Multimedia*, vol. 14, 2012.
- [55] T. Gao and D. Koller, "Discriminative learning of relaxed hierarchy for large-scale visual recognition," in *IEEE International Conference on Computer Vision, ICCV*, 2011.
- [56] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," in *Proceedings of the IEEE*, 1998, pp. 2278–2324.
- [57] M. Zinkevich, M. Weimer, L. Li, and A. J. Smola, "Parallelized stochastic gradient descent," in *Advances in Neural Information Processing Systems 23*, J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, Eds. Curran Associates, Inc., 2010, pp. 2595–2603.

- [58] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *in COMPSTAT*, 2010.
- [59] S. Azadi and S. Sra, "Towards an optimal stochastic alternating direction method of multipliers," in *Proceedings of the 31st International Conference on Machine Learning*, vol. 32. JMLR, 2014, pp. 620–628.
- [60] H. Ouyang, N. He, L. Q. Tran, and A. Gray, "Stochastic alternating direction method of multipliers," in *ICML*, 2013.
- [61] W. Ding, R. Wang, F. Mao, and G. Taylor, "Theano-based large-scale visual recognition with multiple gpus," *arXiv preprint arXiv:1412.2302*, 2014.
- [62] Theano Development Team, "Theano: A Python framework for fast computation of mathematical expressions," *arXiv e-prints*, vol. abs/1605.02688, May 2016.
- [63] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.
- [64] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, "Decaf: A deep convolutional activation feature for generic visual recognition," *CoRR*, vol. abs/1310.1531, 2013.



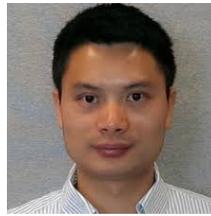
**Tianyi Zhao** received her BS degree in software engineering from Xiamen University in 2015. She is currently pursuing her PhD degree on Computer Science at UNC-Charlotte. Her research interests include semantic image/video classification and retrieval, statistical machine learning, and large-scale visual recognition.



**Baopeng Zhang** received his Ph.D. degree in computer science from Tsinghua University in 2008. He is currently an assistant professor in the School of Computer and Information Technology, Beijing Jiaotong University, China. He was a visiting scholar at UNC-Charlotte from 2015-2016. His research interests include semantic image/video classification and retrieval, statistical machine learning, large-scale semantic data management and analysis, and image privacy protection.



**Wei Zhang** received the BA and MA degrees in economics and the PhD degree in computer science from Fudan University, Shanghai, China, in 2000, 2003, and 2008, respectively. He is currently an associate professor in the School of Computer Science, Fudan University. He is now a visiting scholar at the computer science department in University of North Carolina at Charlotte. His current research interests include machine learning, computer vision, and deep neural network.



**Ning Zhou** received his Ph.D. degree in computer science from UNC-Charlotte in 2013. He is currently a research scientist at Microsoft. His research interests include semantic image/video classification and retrieval, statistical machine learning, large-scale semantic data management and analysis.



**Jun Yu** (M13) received his BEng and PhD from Zhejiang University, Zhejiang, China. He is currently a Professor with the School of Computer Science and Technology, Hangzhou Dianzi University. He was an Associate Professor with School of Information Science and Technology, Xiamen University. From 2009 to 2011, he worked in Singapore Nanyang Technological University. From 2012-2013, he was a visiting researcher in Microsoft Research Asia (MSRA). He was a short-term visiting scholar at UNC-Charlotte. Over the past years, his research interests include multimedia analysis, machine learning, image processing and image privacy protection. He has authored and co-authored more than 50 scientific articles. He has (co-)chaired for several special sessions, invited sessions, and workshops. He served as a program committee member or reviewer top conferences and prestigious journals. He is a Professional Member of the IEEE, ACM and CCF.



**Jianping Fan** is a professor at UNC-Charlotte. He received his MS degree in theory physics from Northwestern University, Xian, China in 1994 and his PhD degree in optical storage and computer science from Shanghai Institute of Optics and Fine Mechanics, Chinese Academy of Sciences, Shanghai, China, in 1997. He was a Postdoc Researcher at Fudan University, Shanghai, China, during 1997-1998. From 1998 to 1999, he was a Researcher with Japan Society of Promotion of Science (JSPS), Department of Information System Engineering, Osaka University, Osaka, Japan. From 1999 to 2001, he was a Postdoc Researcher in the Department of Computer Science, Purdue University, West Lafayette, IN. His research interests include image/video privacy protection, automatic image/video understanding, and large-scale deep learning.