

# Fusing Multiple Hierarchies for Semantic Hierarchical Classification

Shuo Zhao and Quan Zou

**Abstract**—This paper studies the problem of constructing a suitable hierarchy for hierarchical classification. It presents a new method to fuse multiple similarity relatedness between concepts. The method is based on the kernel target alignment technology. We also develop a method to construct a hierarchy for image classification automatically. The hierarchy is constructed based on the previous fused similarity measure. Then, we utilize the structured support vector machine (SVM) for classification with a meaningful hierarchy. Experiments on tow real-world datasets show that hierarchical classification perform better than flat classification, and the structured SVM with the fused classes hierarchy provides a better image classification.

**Index Terms**—Hierarchies construction, hierarchical classification, taxonomies, structural learning.

## I. INTRODUCTION

With the development of information technology, multimedia data grow rapidly. A large amount of data contains thousands classes and varies significantly in semantics. For real-world applications, they are faced the scalability problem which dealing with a huge number of object classes. In practical, humans can easily categorize at least ten thousands of objects and scenes [1]. Humans handle large number of objects through a hierarchical structure built in the semantic space. By means of the hierarchy, humans make fast and meaningful recognition. Several studies [2]-[4] have focused on exploiting hierarchical structure for large-scale classification of text and visual application.

Interpretation of image objects at high level semantic is one of the most critical problems in computer vision. Although many researches have made progress in object recognition, the machine recognition performance is still far from the human intelligence. The process of humans understanding objects happen in high-level semantic space, however most of the machine learning approaches just interpret the object through learning from low-level features. Those approaches can depict the virtual content of the images but they are unable to understand the semantic meanings of the images like humans do. Humans categorize objects in a hierarchical way, which take the similarity between two concepts into count. For example, one may classify a wolf into a dog by mistake, but hardly categorize a wolf as a car. A natural way to incorporate the similarity relatedness into classification is utilize the hierarchical taxonomy [5]-[7].

Manuscript received December 18, 2015; revised February 2, 2016. This work was supported in part by National Natural Science Foundation of China under Grant 61432011.

The authors are with School of Computer Science and Technology, Tianjin University, Tianjin 300350, P. R.China (e-mail: szhao@tju.edu.cn, zouquan@nclab.net).

Several methods have been proposed to construct semantic hierarchies for classification. Many works [2], [5], [8] construct hierarchies based on the well-known large taxonomy—WordNet [9], which group words into sets according their superior-subordinate relations. A semantic hierarchy is constructed based the WordNet in [8]. They extract the relevant subgraph linked all the given concepts from the WordNet. In [2], they propose a large-scale image hierarchy classifier based on the WordNet. The above approaches utilize the high-level conceptual information to build the hierarchy. These conceptual hierarchies contain semantic meaning of the concepts, but they ignore the rich visual information which is also beneficial to classification.

Some approaches [10]-[11] utilize the visual feature information to construct the hierarchies. Bart *et al.* [10] propose a completely unsupervised non-parametric bayesian model to learn a tree hierarchy. Marcin & Cordelia [11] use the  $\chi^2$  distance to compute dissimilarity relatedness between classes and construct a relaxed hierarchy by recursively splitting the class sets until they contain only one class.

Only visual or conceptual information is insufficient to depict the rich contents of images. So some works exploit multiple relatedness information, such as visual feature similarity, tags of images and so on. A graphical model is exploited to automatically construct a “semantivisual” hierarchy which using both visual feature and tags information in [12]. Fan *et al.* [13] construct a hierarchy by adding weighted visual similarity and conceptual similarity. Bannour & Hudelot [14] incorporate three types of similarity information (visual, conceptual, and contextual) to develop a “semantico-visual relatedness of concepts” similarity measure, which is used to construct a faithful hierarchy.

Various classification techniques have been proposed for categorization with a hierarchical structure. In order to improve the efficiency of classification, Gregory & Pietro [15] adopt a complete top-down greedy strategy along the hierarchy, which chooses only the most probable child at each node and ignores other unlikely children until reaching a leaf node. Some other studies [5], [13] also exploit hierarchies to improve the classification performance. In [5], the semantic distance between two nodes in the hierarchical tree is used to define the hierarchical loss, which is used to penalize the misclassifications semantically. Fan *et al.* [13] propose a multi-task learning algorithm that simultaneously learn the correlated classifiers for the sibling nodes sharing the same parent in the hierarchy. Another hierarchical classification approach that exploits the entire tree structure is the popular structured learning framework [6], [16]. In the structured framework [6], the features map function and hierarchical loss function are defined according to the hierarchy, and then learning the hierarchical classifier based

on the structural support vector machines.

Exploiting semantic hierarchies can help categorize objects in a way. However, there are two issues may complicate the usage of hierarchies for classification. First, the conceptual relatedness in the hierarchy is inconsistent with the visual features [17]. Some structure in the hierarchy may damage the classification. For instance, whale and human are semantically similar in the WordNet, but their visual features are dramatically diverse. Shark and whale are fairly semantic distant, but they share some visual features. Second, single semantic hierarchy is insufficient to depict the complex image content. Some previous studies [8], [10], [15] only consider one semantic relatedness. However, in reality, objects have different degrees of relatedness based on different views (e.g., conceptual similarity based on WordNet or visual similarity based on visual features.) Thus, some approaches [13], [14] incorporate multiple similarity relatedness to construct the hierarchy. However, the main drawback of these methods is its similarity measure is just summing the relatedness matrices using weights pre-defined by human.

To address above problems, we present an approach to construct a suitable hierarchy, which utilizes kernel target alignment method [18] to fuse multiple similarity relatedness. We also develop a method to build a hierarchy automatically. For the hierarchical classification with a semantic hierarchy, we translate the hierarchical classification problem into the structured learning framework. We also define a semantic loss function based on the fused similarity measure.

Our main contributions are fusing multiple relatedness measures to construct an appropriate hierarchy, developing a method to build a hierarchy automatically, translating the hierarchical classification problem into the structured learning framework, and defining a novel semantic loss function. We demonstrate our framework on PASCAL VOC and a subset of the Animals with Attributes dataset. The results of a thorough experiment are reported, where the structured SVM with our fused semantic hierarchy provided better performance than flat approaches.

The rest of the paper is organized as follows. The next section introduces the proposed fusing relatedness approach, constructing hierarchy method and structured SVM framework. The experimental results and analysis are given in Section III. Finally, Section IV concludes this paper.

## II. FUSING MULTIPLE SEMANTIC RELATEDNESS

Our goal is to construct a suitable hierarchy for hierarchical classification. We want to fuse multiple semantic relatedness to construct the hierarchy automatically. In order to determine the weight of each semantic relatedness, we adopt the alignment-based techniques [18], which can learning the weight of each kernel in multiple kernel learning (MKL). And then the hierarchy is constructed by a top-down clustering method. In the following, we first introduce the fusing algorithm and the constructing method in Sec. II.A and Sec. II.B, then we describe the structured learning framework in Sec. II.C.

### A. Fusing Multiple Semantic Relatedness

We have a dataset,  $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\} \in X \times Y$ , where

$x_i \in \mathcal{R}^d$  denotes the  $i$ -th example, and  $y_i \in \{1, 2, \dots, C\}$  is its class label, where  $C$  is the number of classes. We want to fuse multiple similarity relatedness, and each similarity relatedness can be expressed as a symmetric matrix  $K \in \mathcal{R}^{C \times C}$ , where  $K(y_i, y_j) \in [0, 1]$ , represent the similarity between classes  $y_i$  and  $y_j$ .

Each symmetric matrix  $K$  can be regarded as a kernel matrix, because the element in kernel matrix represents the dot product of two examples and the dot product can measure the similarity of two examples. Then, we utilize the kernel target alignment (KTA) objective function [18] to learn the weight of each similarity matrix. We define the ideal similarity matrix as the identity matrix  $K_I$ .

For the given  $M$  similarity matrices  $K_1, K_2, \dots, K_M$  and the ideal similarity matrix  $K_I$ , we use Min-Max Normalization to normalize each  $K$  into the same interval. To avoid the kernel scaling problem, each matrix  $K$  is centered by the following equation:

$$K_{cm}(x_i, x_j) = K_m(x_i, x_j) - \frac{1}{C} \sum_{i=1}^C K_m(x_i, x_j) - \frac{1}{C} \sum_{j=1}^C K_m(x_i, x_j) + \frac{1}{C^2} \sum_{i,j=1}^C K_m(x_i, x_j) \quad (1)$$

Following, each  $K_{cm}(x_i, x_j)$  is normalized to have trace equal to one. Then these  $M$  similarity matrices are linearly weighted combined by

$$K_{Similar}(x_i, x_j) = \sum_{m=1}^M w_m K_{cm} \quad (2)$$

where  $w_m$  is the weight of the corresponding  $K_{cm}$  and  $0 \leq w_m \leq 1, \sum_{m=1}^M w_m = 1$ . The weights are learned by maximizing kernel target alignment objective function:

$$\rho(K_{Similar}, K_I) = \frac{K_{Similar} \cdot K_I}{\sqrt{(K_{Similar} \cdot K_{Similar}) \times (K_I \cdot K_I)}} \quad (3)$$

This optimization problem can be efficiently solved by quadratic program. For optimization implementation details, please refer to [18].

The definition of ideal similarity matrix  $K_I$  is to handle the inconsistencies between different relatedness measures. If two classes are similar (or dissimilar) in both the visual and conceptual measure, after summing the linear weighted matrices, the two classes still maintain similar (or dissimilar) due to restricting the sum of weights to one. However, if two classes are inconsistent on visual and conceptual relatedness, the objective function will give a small weight to a large similarity value. For example, the visual similarity between whale and human is smaller than conceptual similarity:  $K_{visual}(whale, human) \ll K_{conceptual}(whale, human)$ . We will obtain a smaller weight for  $K_{conceptual}$  and this can deal with the inconsistencies problem.

### B. Semantic Hierarchy Constructing

After learning the weights of every similarity matrix, we obtain the semantic similarity measure between classes. We

develop a method that construct the semantic hierarchy by top-down recursive partitioning the set of classes until every set only contains one class. At each step, one class set is partitioned into small sets, in which classes are more similar. Each partition procedure can be regarded as a graph cut problem, and it can be solved by standard spectral clustering. We use the Self-Tuning Spectral Clustering [19] to split set. This is a variant of the Spectral Clustering algorithm which can automatically select the number of clusters. Thus, we need not to decide the number of small sets at each partition step. For the affinity matrix in spectral clustering, we use the similarity measure matrix instead of it. By recursively partitioning the classes set, we can obtain a semantic hierarchy. The hierarchy constructed by our method on the Pascal VOC dataset is depicted in Fig. 1, and the hierarchy constructed on the Animal with Attribute dataset is depicted in Fig. 2.

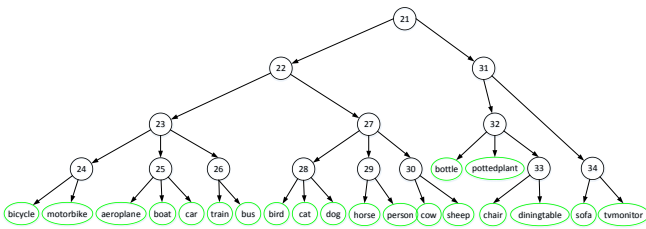


Fig. 1. Semantic hierarchy built on Pascal VOC dataset.

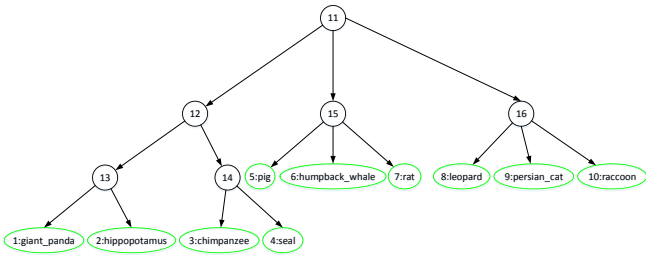


Fig. 2. Semantic hierarchy constructed on Animal with Attribute dataset.

### C. Structural Learning with Taxonomies

Given a hierarchy, we translate the hierarchical classification into the structured SVM learning framework. A hierarchy can be defined as a directed graph  $T=(V,E)$ , where  $V=(v_1, v_2, \dots, v_{|V|})$  and  $Y \in V$  are identified with leaf nodes. The path from root to one leaf node  $y$  is defined as a set of nodes  $\pi(y)$ . The set  $\pi(y)$  can be encoded by a binary vector  $\lambda(y)$ , where the  $i$ -th element is given by

$$\lambda_i(y) = \begin{cases} 1 & \text{if } v_i \in \pi(y) \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

So the class leopard in Fig. 2 is represented by  $\lambda(\text{leopard}) = (0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1)$ .

In structured SVM framework, the discrimination function is defined as

$$f(x_i, y) = \langle w, \phi(x_i, y) \rangle \quad (5)$$

In the hierarchical structure, the mapping function  $\phi(x_i, y)$  is defined as  $\phi(x_i, y) = \lambda(y) \otimes x_i$ , where  $\otimes$  is the tensor product. For the learning objective, we follow the structured

learning framework formulated by Tsochantaridis & Joachims [16]. We formulate the learning objective with a margin rescaling as:

$$\arg \min_{w, \xi} \frac{1}{2} w^T w + \frac{C}{N} \sum_i \xi_i \quad \text{s.t. } \forall i, \xi_i > 0, \quad (6)$$

$$\forall i, \forall \hat{y} \neq y_i : \langle w, \phi(x_i, y_i) - \phi(x_i, \hat{y}) \rangle \geq \Delta(y_i, \hat{y}) - \xi_i$$

where  $C > 0$  is a constant that controls the tradeoff between training error minimization and margin maximization. In this minimization problem (6), the constraint is added to every training instances and each constraint corresponds a slack variable  $\xi_i$ , which is added as an upper bound on the error  $\Delta(y_i, \hat{y})$ . This will make violating a margin constraint with a high  $\Delta(y_i, \hat{y})$  value incurs a more severe penalty. The optimization problem can be solved using the cutting plane algorithm in the SVMStruct software package [20].

The loss function in the structured framework is based on the given ground truth label  $y$  and the predicted label  $\hat{y}$ . In order to classify images more semantically, we take the hierarchical loss into count. In the hierarchical loss criterion, misclassifying an image into a wrong but semantically close class should suffer a smaller loss than misclassifying it into a semantically distant class. Based on this consideration, we define loss function as:

$$\Delta_{\text{fused}}(\hat{y}, y) = 1 - \sum_{m=1}^M w_m K_m(\hat{y}, y) \quad (7)$$

This semantic loss function is based on the fused similarity measure and can describe the meaningful distance between the true label  $y$  and the prediction  $\hat{y}$ .

## III. EXPERIMENTS

We perform an empirical study on two real-world datasets, Pascal VOC [21] and Animals with Attributes [22].

### A. PASCAL VOC

We perform experiments on the Pascal VOC2010 dataset, which contains 20 classes and a total of 11,321 images. For the features, we adopt a bag of words (BoW) model. The BoW feature vector is built as following: compute dense scale-invariant feature transform (SIFT) descriptors, generate codebook and encoding SIFT. The dense SIFT features are computed using the VLFeat package [23]. Then, the dense features are further processed to form a visual codebook of  $D=1000$  visual words using K-means clustering. Every dense feature in one image is assigned to one visual word in the codebook through a KD-Tree, and the image is represented by a histogram of  $D$  visual words. We also adopt the explicit feature map approach [24], that enable linear SVM obtain comparable performance to the nonlinear SVMs with implicit kernels. In our setups, the  $\chi^2$  kernel is adopted with approximation order  $N=1$ . Then, one image is represented as a BoW vector with 1000 dimensions.

For the multiple similarity relatedness, we computed three similarity relatedness measures as follows [14]: visual similarity, which computes the similarity between centers of different classes; conceptual similarity, which represents the

distance of paths connecting two concepts in WordNet; and the co-occurrence probability between each pair of concepts. We performed classification on Pascal VOC2010 dataset using the structured SVM. We trained on the divided training set and tested on the validation set supplied by the dataset. The classification performance was evaluated using the

average precision (AP) score, a standard evaluation criterion supplied by the PASCAL challenge. The AP score represent the area under the precision/recall curve, and a higher value indicates a better performance. The mean average precision (mAP), which computes the mean AP score for the 20 classes, is reported.

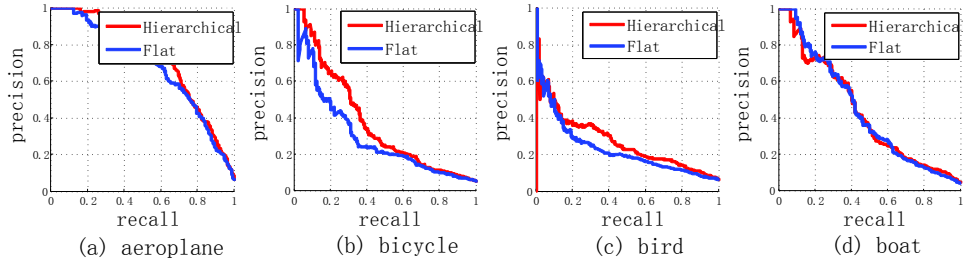


Fig. 3. Precision/recall curve for Pascal VOC2010. The fused hierarchical method performs better than the flat method on most classes. We only listed the first four classes due to space constraints. The other classes are similar with the same relationship between two methods.

We firstly compared the hierarchical classification with the flat classification, which only adopts the joint feature representation but ignores the hierarchy structure. Results are shown in Fig. 3. The fused hierarchical classification performs better than the flat classification on most classes. This shows that exploiting hierarchies can enhance the recognition performance. In order to evaluate our fused hierarchy, we compared our hierarchy (SSVM+Fused) with the hierarchy (SSVM+Hichem) built in [14]. For further comparison, we additionally reported the result on three other methods: random forest (RF), flat structured multiclass SVM (Flat-SVM), and linear-SVM. The results are shown in Fig. 4. The structure learning framework with our fused hierarchy shows the best performance than the other methods. Our fused hierarchy performs an improvement of 1.97% than the hierarchy in [14] using the same classifier. The result demonstrates that our fused hierarchy can more accurately depict the semantic relatedness of the categories.

computed the Euclidean distance between real-valued attributes vectors supplied in the dataset for the training images.

We split the images into 100/100/100 images per class for training/validation/testing, and generated five such random splits. We reported the average recognition accuracy and standard errors for a 95% confidence interval. Our method is compared with several popular multi-class classification algorithms, which include random forest (RF), flat structured SVM (Flat-SVM), and linear-SVM. We also compared our method with semantic kernels forests (SKF) method [17], which exploits multiple hierarchies form multiple views. The results are summarized in Fig. 5. Our classification framework produces the best results than others, which include flat and hierarchical methods. The results show that the hierarchical method outperforms the flat algorithms. Our method is also comparable with the SKF method, which exploits multiple kernel learning and nolinear kernels.

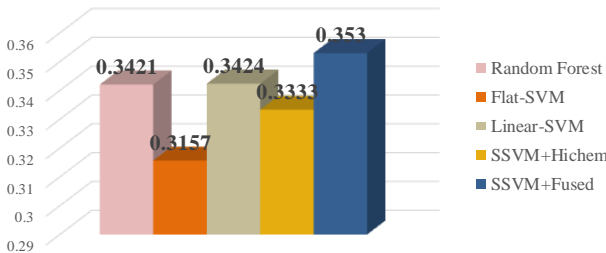


Fig. 4. Performance of different methods on Pascal VOC2010.

**B. Animals with Attribute**

For the Animals with Attribute dataset, we used a subset that contained ten classes in [17], and a total of 6,180 images. To obtain a rich image features, we made use of the deep convolutional activation features (DeCAF) supplied by the dataset. In order to comparing with the semantic kernels forests (SKF) method [17], these features are reduced to 100 dimensionality using principal component analysis (PCA).

For the multiple similarity relatedness, we computed four similarity relatedness measures: visual similarity and conceptual similarity following PASCAL VOC, and two other relatedness, appearance similarity, and habitat similarity. The appearance similarity and habitat similarity were computed followed the method in [17]. The method

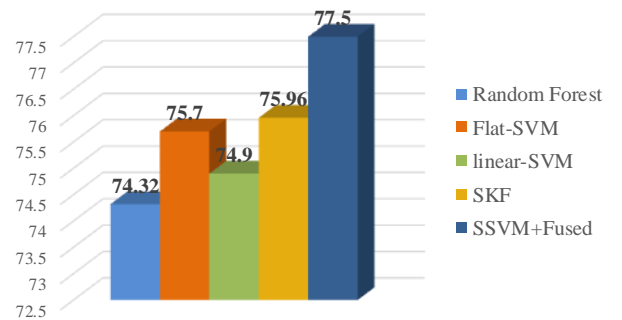


Fig. 5. Performance of different methods on Animals with Attribute.

**C. Comparison with Different Loss Function**

Base on the intuition that misclassifying an image should incurs different semantically penalty, different loss functions can be defined based on the hierarchy. We compared four loss functions, which include the standard 0/1 loss hierarchy-based loss  $\Delta_{0/1}$  [25], hierarchy-based loss function  $\Delta_h(\hat{y}, y) = \sum_i |\lambda_i(\hat{y}) - \lambda_i(y)|$ , weighted hierarchical difference(WHD) loss [26], and our semantic loss. The weighted hierarchical difference is defined as:

$$\Delta_{WHD}(\hat{y}, y) = \sum_i |\Psi(\lambda_i(\hat{y})) - \Psi(\lambda_i(y))| \tag{8}$$

which penalizes more severely, when the misclassification occurs higher level of the hierarchy.  $\Psi(\lambda_i(y))$  is defined as a weighting function, which divides each element of the binary vector  $\lambda(y)$  by its level. For example, in the hierarchy of Fig. 2, the  $\Psi(\lambda(\text{leopard}))$  of class 8:leopard is defined as  $\Psi(\lambda(\text{leopard})) = (0, 0, 0, 0, 0, 0, 0, 1/3, 0, 0, 1/2, 0, 0, 0, 0, 1)$ . We reported the performances using different loss functions on PASCAL VOC and AWA datasets. The results are list in Table I. Our semantic loss function has the best performance. These results illustrate that the fused similarity measure is more meaningful and precisely describe the relatedness between classes.

TABLE I: PERFORMANCES WITH VARIOUS LOSS FUNCTIONS

Dataset	$\Delta_{0/1}$	$\Delta_h$	$\Delta_{WHD}$	$\Delta_{fused}$
AWA-10	76.64 $\pm$ 1.31	75.70 $\pm$ 1.01	75.71 $\pm$ 1.28	<b>77.50<math>\pm</math>1.13</b>
VOC2010	0.3506	0.3101	0.2926	<b>0.3530</b>

#### IV. CONCLUSIONS

We presented an approach that can fuse multiple relatedness measures to construct a class hierarchy for hierarchical classification. The experimental results showed that the fused semantic measure can depict the relationship between different concepts, and the proposed method improved object recognition accuracy. In our future work, we plan to explore evaluation criterions to produce more semantic prediction for image classification.

#### REFERENCES

[1] I. Biederman, "Recognition-by-components: A theory of human image understanding," *Psychological Review*, vol. 94, p. 115, 1987.

[2] B. Zhao, F. Li, and E. P. Xing, "Large-scale category structure aware image categorization," *Advances in Neural Information Processing Systems*, pp. 1251-1259, 2011.

[3] S. Bengio, J. Weston, and D. Grangier, "Label embedding trees for large multi-class tasks," *Advances in Neural Information Processing Systems*, pp. 163-171, 2010.

[4] R. Fergus, H. Bernal, Y. Weiss, and A. Torralba, "Semantic label sharing for learning with many categories," *Computer Vision—ECCV*, pp. 762-775, 2010.

[5] J. Deng, A. C. Berg, K. Li *et al.*, "What does classifying more than 10,000 image categories tell us?" *Computer Vision-ECCV*, vol. 6315, pp. 71-84, 2010.

[6] L. Cai and T. Hofmann, "Hierarchical document categorization with support vector machines," in *Proc. the Thirteenth ACM International Conference on Information and Knowledge Management*, 2004, pp. 78-87.

[7] J. Y. Chang and K. M. Lee, "Large margin learning of hierarchical semantic similarity for image classification," *Computer Vision and Image Understanding*, vol. 132, pp. 3-11, 2015.

[8] M. Marszałek and C. Schmid, "Semantic hierarchies for visual object recognition," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1-7, 2007.

[9] G. Miller and C. Fellbaum, "Wordnet: An electronic lexical database," MIT Press Cambridge, 1998.

[10] E. Bart, I. Porteous, P. Perona, and M. Welling, "Unsupervised learning of visual taxonomies," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1-8.

[11] M. Marszałek and C. Schmid, "Constructing category hierarchies for visual recognition," *Computer Vision—ECCV*, 2008, pp. 479-491.

[12] L.-J. Li, C. Wang, Y. Lim, D. M. Blei *et al.*, "Building and using a semantivisual image hierarchy," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2010, pp. 3336-3343.

[13] J. Fan, Y. Gao, and H. Luo, "Integrating concept ontology and multitask learning to achieve more effective classifier training for multilevel image annotation," *IEEE Transactions on Image Processing*, vol. 17, pp. 407-426, Mar. 2008.

[14] H. Bannour and C. Hudelot, "Building semantic hierarchies faithful to image semantics," *Advances in Multimedia Modeling*, vol. 7131, pp. 4-15, 2012.

[15] G. Griffin and P. Perona, "Learning and using taxonomies for fast visual categorization," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1-8.

[16] I. Tschantzaris, T. Joachims, T. Hofmann, and Y. Altun, "Large margin methods for structured and interdependent output variables," *Journal of Machine Learning Research*, 2005, pp. 1453-1484.

[17] S. J. Hwang, K. Grauman, and F. Sha, "Semantic kernel forests from multiple taxonomies," *Advances in Neural Information Processing Systems*, 2012, pp. 1718-1726.

[18] C. Cortes, M. Mohri, and A. Rostamizadeh, "Two-stage learning kernel algorithms," in *Proc. the 27th International Conference on Machine Learning*, 2010, pp. 239-246.

[19] L. Zelnik-Manor and P. Perona, "Self-tuning spectral clustering," *Advances in Neural Information Processing Systems*, pp. 1601-1608, 2004.

[20] T. Joachims. (2008). Support vector machine for complex outputs. [Online]. Available: [http://www.cs.cornell.edu/people/tj/svm\\_light/svm\\_struct.html](http://www.cs.cornell.edu/people/tj/svm_light/svm_struct.html)

[21] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International Journal of Computer Vision*, vol. 88, pp. 303-338, 2010.

[22] C. H. Lampert, H. Nickisch, and S. Harmeling, "Learning to detect unseen object classes by between-class attribute transfer," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 951-958.

[23] A. Vedaldi and B. Fulkerson, "VLFeat: An open and portable library of computer vision algorithms," in *Proc. the International Conference on Multimedia*, 2010, pp. 1469-1472.

[24] A. Vedaldi and A. Zisserman, "Efficient additive kernels via explicit feature maps," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, pp. 480-492, 2012.

[25] A. Binder, K.-R. Müller, and M. Kawanabe, "On taxonomies for multi-class image categorization," *International Journal of Computer Vision*, vol. 99, pp. 281-301, Sep. 2012.

[26] N. Nourani-Vatani, R. López-Sastre, and S. Williams, "Structured output prediction with hierarchical loss functions for seafloor imagery taxonomic categorization."



**Shuo Zhao** is currently a master candidate with the School of Computer Science and Technology in Tianjin University. His research interests include computer vision and machine learning.



**Zou Quan** is a professor of computer science at Tianjin University. He received his Ph.D. from Harbin Institute of Technology, P.R.China in 2009. From 2009 to 2015, he is an assistant and associate professor in Xiamen University, P.R.China. His research is in the areas of bioinformatics, machine learning and parallel computing.