

Convolutional Fisher Kernels for RGB-D Object Recognition

Yanhua Cheng^{1,2,4}, Rui Cai⁵, Xin Zhao^{1,2,4}, and Kaiqi Huang^{1,2,3,4}

¹Center for Research on Intelligent Perception and Computing

²National Laboratory of Pattern Recognition

³CAS Center for Excellence in Brain Science and Intelligence Technology

⁴Institute of Automation, Chinese Academy of Sciences

⁵Microsoft Research

Abstract

This paper studies the problem of improving object recognition using the novel RGB-D data. To address the problem, a new convolutional Fisher Kernels (CFK) method is proposed to represent RGB-D objects powerfully yet efficiently. The core idea of our approach is to integrate the both advantages of the convolutional neural networks (CNN) and Fisher Kernel encoding (FK): CNN model is flexible to adapt to new data sources, but requires for large amounts of training data with significant computational resources for good generalization; In comparison, FK encoding is able to represent objects powerfully and efficiently with small training data, however, its success highly depends on the well-designed SIFT features in literature, which may not be suitable for the new depth data. CFK can be interpreted as a two-layer feature learning structure to bridge the two models. The first layer employs a single-layer CNN to learn low-level translationally invariant features for both RGB and depth data efficiently. The second layer aggregates the convolutional responses by FK encoding. Here 2D and 3D spatial pyramids are applied to further improve the Fisher vector representation of each modality. Experiments on RGB-D object recognition benchmarks demonstrate that our approach can achieve the state-of-the-art results.

1. Introduction

The past several years have witnessed the rapidly increasing popularity of object recognition by fusing the RGB and depth (RGB-D) data [19, 8, 3, 5, 20, 4, 7, 30, 26, 27, 27, 10]. Such a booming application of the novel RGB-D data mainly attributes to recent depth cameras, such as Kinect, which are able to record high quality frames of both color and depth information synchronously. Another remarkable

trend is that these advanced depth cameras are being integrated into mobile devices like Google Tango [1] and Microsoft HoloLens [2], which further promote the researches of RGB-D object recognition.

Although the captured RGB-D data provides rich multi-modal information to depict an object, such as color, texture, appearance (RGB modality) as well as shape and geometry information (depth modality), how to effectively represent each modality and combine the both to improve object recognition remains an open problem. Much progress has been made in the past few years, from bag-of-words model with handcrafted features [19, 8, 26, 33] and efficient match kernels with kernel features [5] to feature learning approaches [4, 7, 30].

Motivated by the great success of Fisher Kernel encoding (FK [24, 25]) and deep convolutional neural networks (CNN [18]) in visual object recognition such as Pascal VOC Challenge [14] and ImageNet Challenge [13], this paper proposes a new feature learning method for RGB-D object recognition, termed convolutional Fisher Kernels (CFK). CFK is developed to learn features from the raw RGB and depth data powerfully and efficiently by combining the unique advantages of FK and CNN model. Towards FK, it is a generic framework to combine the benefits of generative and discriminative approaches. In the context of image classification, it has yielded clearly superior results on a variety of image benchmarks than the popular bag-of-words (BoW) model based on other feature encoding methods such as LLC [32], super vector encoding [34], etc. (see [9, 16] for a survey). Another major advantage of FK is that a very small codebook size and linear classifiers are sufficient for these impressive results, i.e., computing FK feature vectors and learning classifiers are very efficient in practice. However, in literature, the success of FK for image recognition [24, 25] is highly dependent on the well-designed SIFT features, which are not suitable for the novel depth data as

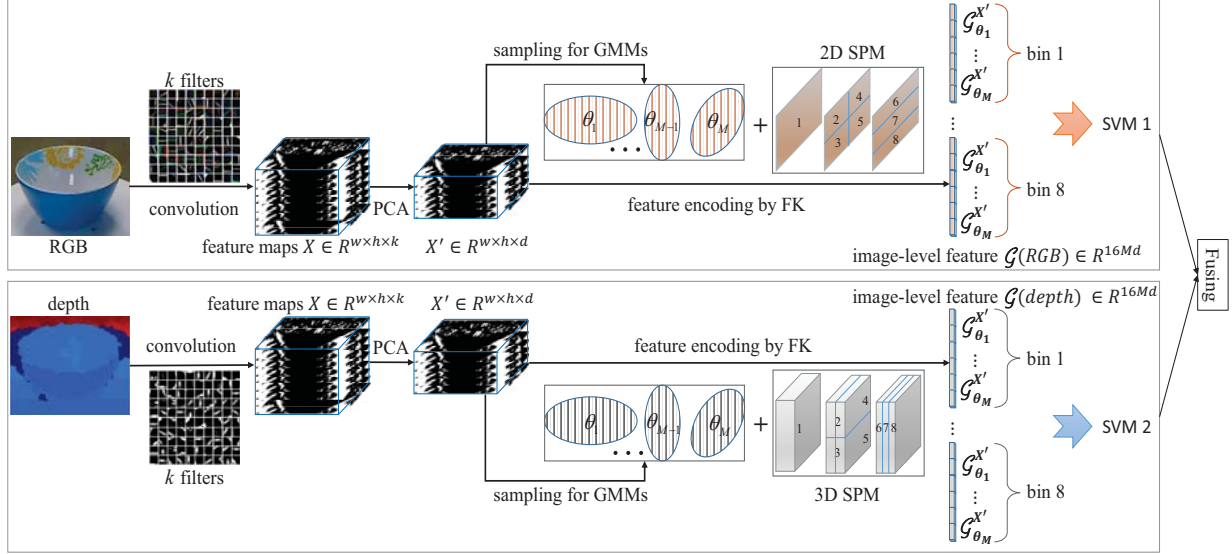


Figure 1. Overview of the CFK framework for RGB-D object recognition. Best viewed in color.

shown in [5, 8]. To get rid of this drawback, a *single-layer CNN model* (It can be extended to multi-layer architecture easily, however, we empirically find that a single layer is sufficient for CFK to achieve the state of the art result on RGB-D object recognition benchmarks [19, 8] whilst being efficient.) is employed in this paper to learn low-level translationally invariant features from the raw RGB and depth data, respectively. Finally, in analogy to spatial pyramid matching in BoW model [21], *2D and 3D spatial pyramids* are exploited separately for the two modalities to further improve the performance of the FK representation. To the best of our knowledge, this is the first work to combine convolutional neural networks and Fisher Kernel encoding for RGB-D object recognition. Experimental results demonstrate that our approach can significantly outperform the existing state-of-the-art methods.

The rest of this paper is organized as follows. Section 2 briefly reviews related work. Section 3 introduces the proposed CFK approach for RGB-D object recognition. Extensive experimental results are reported in Section 4 and conclusion is drawn in Section 5.

2. Related Work

In this section, we discuss the relevant prior work with focus on (1) methods for RGB-D object recognition and (2) approaches that attempt to bridge FK and CNN model.

Methods for RGB-D Object Recognition. Most of research efforts focused on feature extraction. Lai *et al.* [19] extracted handcrafted features such as SIFT [23], texon histograms [22] and spin images [17] over the color and depth frames individually, followed by a concatenation of all these features to depict each object. In view of the specificity of

the depth information, quantized 3D SURF local descriptors with selective 3D spatial pyramids [26] and depth kernel features [5] were designed to represent the depth cue more effectively. Another line of work [4, 7, 30] exploited very successful machine learning methods to learn powerful features from the raw RGB-D data, and obtained very promising results for object recognition. Instead of feature extraction, Lai *et al.* [20] paid attention to RGB and depth fusion by defining a view-to-object distance via sparse distance metric learning, and Cheng *et al.* [10, 11] proposed a semi-supervised framework to make use of both the labeled and unlabeled RGB-D data, which benefited a lot from the complementarity between the RGB and depth cues. In this paper, the proposed CFK approach belongs to the family of feature learning methods.

Deep Fisher Networks. There have been attempts to bridge FK and CNN for visual object recognition. Simonyan *et al.* [29] stacked FK encoding recursively in multiple layers analogous to the deep architecture of CNN model. It achieved competitive results with CNN at a smaller computational cost. Very recently, Sydorov *et al.* [31] proposed to train Fisher Kernel SVMs in a deep way, of which the classifier parameters and GMM parameters are learned and optimized jointly from training data. Indeed, both of the two work tried to construct the deep Fisher networks like CNN model to represent objects more powerfully. However, they still depended on the well-designed SIFT features to obtain the feature responses of the first layer. SIFT feature are popular to characterize the contents in color or grayscale images, but are not suitable to new data sources like depth frames. To address the problem, the proposed CFK in this paper is a more explicit way to combine the CNN model and the FK encoding, which can learn powerful features from

the raw RGB and depth data, whilst being efficient.

3. Convolutional Fisher Kernels

3.1. Overview

Fig. 1 illustrates the overview of the CFK framework for RGB-D object recognition. Given the RGB or depth modality of an object, we first learn its convolutional feature maps via the pretrained filters. Then PCA decorrelation is applied and lower-dimensional convolutional responses are obtained. After that, Gaussian mixture models are utilized to model the distribution of these convolutional feature vectors, and 2D or 3D spatial pyramids are introduced to consider the rough geometry structure information of an object. Finally, Fisher Kernel encoding is exploited to compute the gradient of the log-likelihood of the distribution probability density function and generate the image-level feature vector. A linear SVM classifier is trained for each modality, and the combined scores are used to predict the category label of a new object I_t :

$$\text{category}(I_t) = \underset{c_i \in \mathcal{C}}{\operatorname{argmax}} \gamma S_{RGB}^{c_i} + (1 - \gamma) S_{depth}^{c_i}, \quad (1)$$

where $0 \leq \gamma \leq 1$ is the coefficient for score fusion and \mathcal{C} denotes the label set of all the categories. Details of the process are given below.

3.2. Single-Layer Convolutional Neural Networks

Derived from the deep learning structures, a single CNN layer with the state-of-the-art preprocessing [12] is demonstrated to be powerful yet efficient to learn low-level features from the raw data. This paper applies it for RGB and depth modality, respectively. The procedure of single-layer CNN mainly consists of two steps: (1) pretraining the filters by k-means clustering over the randomly sampled and preprocessed image patches; and (2) convolving the filters over the whole image to generate the feature responses.

Pretraining The Filters. For RGB or depth modality, we first sample a set of squared sub-patches from the training image set. According to [12], each sub-patch is normalized by subtracting the mean and then divided by the standard deviation of its elements. In addition, ZCA whitening is performed to de-correlate pixels and remove redundant features from raw images. Note that these preprocessing steps are crucial, especially for the depth modality, since the absolute depth values can be confused to depict the geometry and shape cues of an object. Finally, k-means clustering is used to learn the filters over these preprocessed sub-patches. Fig. 2 shows the resulted RGB and depth filters. We can see that both RGB and depth filters are sensitive to edge structures of objects, while RGB filters are able to capture additional color cues.

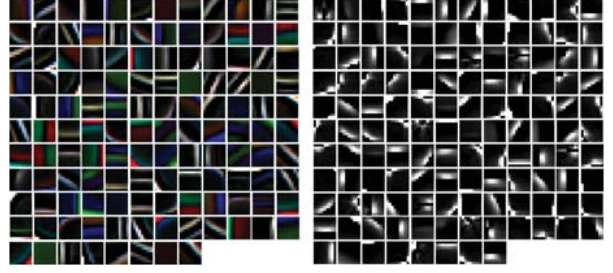


Figure 2. The RGB filters (left) and the depth filters (right) learned by k-means clustering on the Washington RGB-D object dataset [19]. Best viewed in color.

Learning Convolutional Feature Maps. For each modality of each object, we convolve the corresponding filters over the whole image and obtain the convolutional feature maps $X(RGB), X(depth) \in \mathbb{R}^{w \times h \times k}$, where $w \times h$ denotes the size of each feature map, and k is the number of feature maps (equal to the number of the filters). It is worth noting that both $X(RGB)$ and $X(depth)$ are post-processed by rectification with absolute values and local normalization for the convolutional feature responses, which is empirically shown to be important for recognition [12]. Finally, PCA reduction is performed to de-correlate convolutional feature maps of each modality, and generates the reduced feature responses $X'(RGB), X'(depth) \in \mathbb{R}^{w \times h \times d}$, where $d < k$.

3.3. Fisher Kernel Encoding

FK is first introduced by [24] to visual image classification, and then significantly improved by [25]. This paper applies it to new RGB-D object recognition. Taking the convolutional feature responses X' of each modality of each object as a set of input feature vectors $\chi = \{x_1, x_2, \dots, x_N\} \in \mathbb{R}^{d \times N}$ ($N = w \times h$), the target of CFK is to construct the image-level feature vector $\mathcal{G} \in \mathbb{R}^D$ by characterizing χ via Gaussian mixture models. The above process is called Fisher Kernel encoding, and note that it is separate for RGB and depth modality.

Fisher Kernel encoding mainly consists of three modules. First, Gaussian mixture models (GMMs) are learned to describe the probability density distribution of the convolutional feature vectors. The parameters of GMMs are denoted as $\theta = \{(w_m, \mu_m, \Sigma_m)\}_{m=1}^{m=M}$, where w_m , μ_m and Σ_m mean the weight, the mean vector and the covariance matrix of the m th Gaussian model. According to [25], we train the GMMs by using the Maximum Likelihood criterion via a standard Expectation-Maximization algorithm. Second, assume that the convolutional feature vectors χ extracted from an object (RGB or depth modality) are independent each other, then the object can be described by the

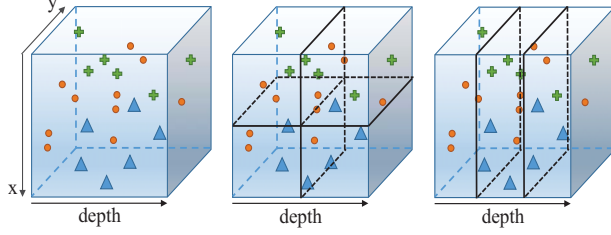


Figure 3. 3D spatial pyramids with Level 0, Level 1, and Level 2 (from left to right).

following gradient vector:

$$\mathcal{L}(\chi|\theta) = \frac{1}{N} \sum_{n=1}^{n=N} \frac{\partial \log p(x_n|\theta)}{\partial \theta}, \quad (2)$$

where $p(x_n|\theta)$ is the probability density function of GMMs. Using the Fisher information matrix [24]

$$F_\theta = E_{x \sim \theta} \left[\frac{\partial \log p(x|\theta)}{\partial \theta} \times \left(\frac{\partial \log p(x|\theta)}{\partial \theta} \right)' \right], \quad (3)$$

the Fisher vector is represented as the normalized gradient vector

$$\mathcal{G}_\theta^\chi = F_\theta^{-1/2} \mathcal{L}(\chi|\theta). \quad (4)$$

Following [24], the diagonal closed-form approximation is used for the Fisher information matrix, and the above Fisher vector can be computed as follows:

$$\begin{aligned} \mathcal{G}_\theta^\chi &= [\mathcal{G}_{\mu,i}^\chi; \mathcal{G}_{\Sigma,i}^\chi], i = 1, \dots, M, \\ \mathcal{G}_{\mu,i}^\chi &= \frac{1}{N\sqrt{w_i}} \sum_{n=1}^N \varphi_n(i) \Sigma_i^{-1/2} (x_n - \mu_n), \\ \mathcal{G}_{\Sigma,i}^\chi &= \frac{1}{N\sqrt{2w_i}} \sum_{n=1}^N \varphi_n(i) [\Sigma_i^{-1} (x_n - \mu_n)^2 - 1], \\ \varphi_n(i) &= \frac{w_i p(x_n|\theta_i)}{\sum_{j=1}^M w_j p(x_n|\theta_j)}. \end{aligned} \quad (5)$$

Here $\varphi_n(i)$ denotes the soft voting of descriptor x_n to Gaussian i . According to [25], the gradient $\mathcal{G}_{w,i}^\chi$ to w is discarded as it brings little additional information. Therefore, the Fisher vector for the holistic image is constructed by the concatenation of $\mathcal{G}_{\mu,i}^\chi$ and $\mathcal{G}_{\Sigma,i}^\chi$ for all the Gaussians. The generated Fisher vector is $\mathcal{G}_\theta^\chi \in \mathbb{R}^{2Md}$ (without spatial pyramids). Finally, 2D and 3D spatial pyramids are applied for RGB and depth modality individually, which are to take into account the rough spatial geometry information of the object for FK encoding. Details are given in the next section.

3.4. 2D/3D Spatial Pyramids

Lazebnik *et al.* [21] demonstrated that the weak spatial geometry information by 2D spatial pyramids can significantly improve the performance of BoW models for visual object/scene recognition. Similarly, Perronnin *et al.* [25]

applied 2D spatial pyramids to improve FK encoding for visual data representation. In analogy to the 2D spatial pyramids, 3D spatial pyramids are utilized for the depth modality in our approach.

As illustrated in Fig. 3, the depth frame is first projected to the three-dimensional space (aka 3D point clouds), and then a three-level spatial pyramid $\{1 \times 1 \times 1, 1 \times 2 \times 2, 1 \times 1 \times 3\}$ are designed. Instead of subdividing the depth frame uniformly, we fine-tune the parting plane to make sure that each sub-region contains the approximately same number of depth pixels. It is necessary since the uniform subdivision can suffer from the problem that some subregions may have few depth pixels and cause noisy representation. With the rectified 3D spatial pyramid, we extract a robust Fisher vector for each sub-region as Section 3.3, and concatenate all the Fisher vectors to represent the depth cue, i.e., with spatial pyramids, the dimension of the image-level feature representation \mathcal{G}_θ^χ is $16Md$ (There are 8 bins in total for the three-level spatial pyramid in Fig. 3). It is worth noting that a similar yet more complex 3D spatial pyramid is also applied to design a new handcrafted depth feature in [26]. Towards the 2D spatial pyramid for the RGB modality, this paper inherits the three-level partitions $\{1 \times 1, 2 \times 2, 1 \times 3\}$, which are adopted by [21, 25].

Following [25], the Fisher vector of each bin is post-processed individually with a two-step normalization: power normalization and L2 normalization:

$$\begin{aligned} f^{(1)} &= \text{sign}(\mathcal{G}_\theta^\chi) |\mathcal{G}_\theta^\chi|^\alpha, \\ f^{(2)} &= f^{(1)} / \sqrt{\|f^{(1)}\|_2^2}, \end{aligned} \quad (6)$$

where $0 \leq \alpha \leq 1$. Then the normalized Fisher vector of each bin are concatenated $[f_{bin1}^{(2)}; f_{bin2}^{(2)}; \dots; f_{bin8}^{(2)}]$ to represent the holistic object. Such a normalization is crucial for object recognition with a linear classifier.

4. Experiments

This paper evaluates CFK for RGB-D object recognition on the challenging Washington RGB-D object dataset [19] and the 2D3D dataset [8], with comparison to the existing state-of-the-art methods.

4.1. Experimental Setup

CFK Feature Representation. Towards the single-layer CNN, the size of each modality of each object is first scaled to 148×148 , and then 128 filters with size $9 \times 9 \times R$ (R denotes the number of channels for the corresponding modality) are pretrained on these scaled images. The resulted convolutional feature response for each modality are $140 \times 140 \times 128$ -dimensional. After PCA reduction, we generate $140 \times 140 \times 80$ -dimensional feature response.

Methods	Depth	RGB	Combine
Linear SVM [19]	53.1 ± 1.7	74.3 ± 3.3	81.9 ± 2.8
Kernel SVM [19]	64.7 ± 2.2	74.5 ± 3.1	83.8 ± 3.5
Random Forest [19]	66.8 ± 2.5	74.7 ± 3.6	79.6 ± 4.0
IDL [20]	70.2 ± 2.0	78.6 ± 3.1	85.4 ± 3.2
3D SPMK [26]	67.8	—	—
KDES [5]	78.8 ± 2.7	77.7 ± 1.9	86.2 ± 2.1
CKM [4]	—	—	86.4 ± 2.3
HMP [6]	70.3 ± 2.2	74.7 ± 2.5	82.1 ± 3.3
SP-HMP [7]	81.2 ± 2.3	82.4 ± 3.1	87.5 ± 2.9
CNN-RNN [30]	78.9 ± 3.8	80.8 ± 4.2	86.8 ± 3.3
CNN-RNN+CT [10]	77.7 ± 1.4	81.8 ± 1.9	87.2 ± 1.1
CNN-SPM-RNN+CT [11]	83.6 ± 2.3	85.2 ± 1.2	90.7 ± 1.1
CNN-only baseline	78.1 ± 1.3	82.7 ± 1.2	87.5 ± 1.1
CFK	85.8 ± 2.3	86.8 ± 2.2	91.2 ± 1.5

Table 1. Comparison of recent results on the Washington RGB-D object database.

For the large-scale Washington RGB-D dataset, we train $M = 256$ Gaussians based on the reduced convolutional feature vectors, while $M = 128$ for the relatively small-scale 2D3D dataset. Towards the two-step normalization, we fix $\alpha = 0.5$ as [25].

Linear SVM Classifiers. For each modality, a linear SVM classifier with a hinge loss is trained based on the training set. We learn the classifier by the primal formulation and a Stochastic Gradient Descent algorithm [28]. Similar to [7], another two modalities, the grayscale image (derived from the RGB modality) and the surface normals (derived from the depth modality) are used in this paper. Given an object, the score S_{RGB} of the RGB modality is actually averaged by the RGB and grayscale scores, and the score S_{depth} is averaged by the depth and normal scores. The coefficient γ is also simply set to 0.5 in Eqn. 1.

4.2. Washington RGB-D Dataset

Dataset. The Washington RGB-D object dataset [19] is a large-scale and multi-view object dataset captured by the Microsoft Kinect. It collects 300 household objects, grouped into 51 categories. Each object instance is imaged from 3 vertical angles as well as multiple horizontal angles, resulting roughly 600 images per instance. There are a total of around 207,920 color and depth images, which are subsampled every 5th frame from each instance and give around 41,877 images for category recognition.

Following the same setting of the work [19], we utilize the provided 10 trials to average the accuracies. For each train/test trial, one object instance is randomly selected from each category for testing, and the remaining object instances are for training.

Results. Table 1 shows the comparison of the recog-

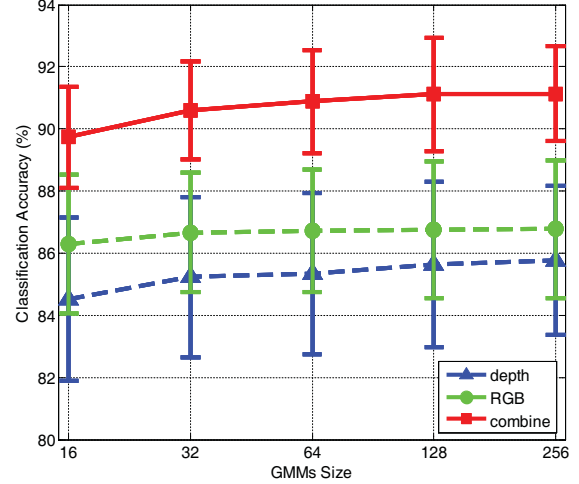


Figure 4. Effect of the GMMs size (M) to the performance of CFK on the Washington RGB-D dataset (2D and 3D spatial pyramids are applied here).

tion results on the Washington RGB-D object dataset. Instead of using handcrafted features, SP-HMP method [7] designed a two-layer hierarchical matching pursuit via sparse coding to learn powerful feature from the raw RGB-D data, which achieved sizable performance gains. CNN-RNN [30] proposed a deeper feature learning model by combining the convolutional neural networks and recursive neural networks. However, it behaved a litter worse than SP-HMP, for which the main reason is probable that the simplified fixed-tree RNNs used in [30] can limit the learning ability of the CNN-RNN model a lot. In addition, we show the result of the popular AlexNet [18] as the CNN-only baseline, which exploits RGB-D data following the method [15]. The CNN-only baseline may suffer from overfitting, as most test instances are unseen in the training set for this challenging dataset. The work [10, 11] proposed a co-training method to fuse RGB and depth modality effectively and achieved the state of the art together with CNN-SPM-RNN features. Our CFK approach can clearly outperform all these methods over each modality as well as the both. Note that all the results of CFK is based on a linear SVM classifier, which is efficient for practice.

Parameter Analysis. The performance of the proposed CFK approach is closely related to the size M of the GMMs as well as the spatial pyramid matching. Each of them is analysed individually by keeping the other one the same with the default setting (i.e., $M = 256$ with 2D/3D spatial pyramids).

(1) The effect of M . As shown in Fig. 4, the recognition accuracy of CFK can rise gradually as M increases. This is because that more Gaussian mixture models can model the distribution of the convolutional feature vectors better, and then more powerful feature representation can be generated by the FK encoding. Note that even a small size, e.g.

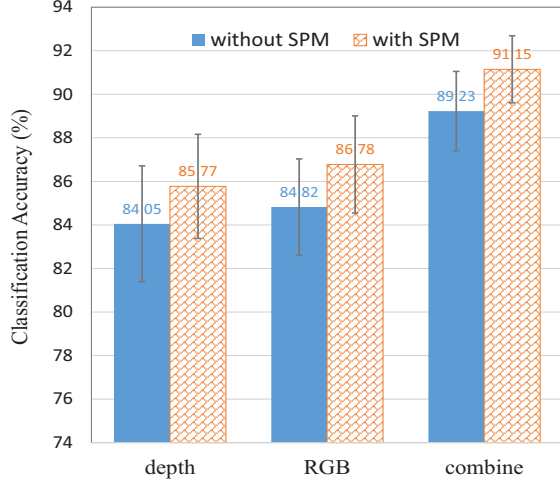


Figure 5. Effect of spatial pyramid matching to the performance of CFK on the Washington RGB-D dataset ($M = 256$).

$M = 16$, can guarantee a superior result for CFK to the previous methods. When $M > 128$, it converges to a stable state with a high accuracy, since very limited additional information can be further attained to the added Gaussians. In this paper, we fix $M = 256$ for the Washington RGB-D dataset.

(2) The effect of the 2D/3D spatial pyramid matching. Fig. 5 compares the recognition accuracies of CFK with and without SPM. For the depth modality, the introduced 3D spatial pyramid can yield a 1.72% improvement. Meanwhile, the 2D spatial pyramid can improve RGB-based object recognition by 1.96%. The final combined result with spatial pyramids is 91.15% accuracy, superior to the one without SPM (89.23% accuracy). The results demonstrate that the rough spatial geometry information (mainly the location of object parts) can be useful for object representation.

Error Analysis. Fig. 6 shows the confusion matrix of CFK over the 51 testing categories of the Washington RGB-D dataset. We can observe that most categories can be correctly classified, which demonstrate the effectiveness of our approach as well. However, there are some easily confused categories such as *pitchers*, *balls*, *mushrooms* and *peaches*, which can have nearly the same shape, appearance or even both with other categories from the certain viewing angles. Some examples are shown in Fig. 7.

4.3. 2D3D Dataset

Dataset. The 2D3D object dataset [8] collects objects from typical household and office environments. It consists of 156 object instances organized into 14 categories. Each object instance is recorded every 10° around the vertical axis on a turntable, yielding 36 views per instance. There are a total of 5,616 RGB-D images. Following the same setting of [8] for category recognition, we first sample 18 views of

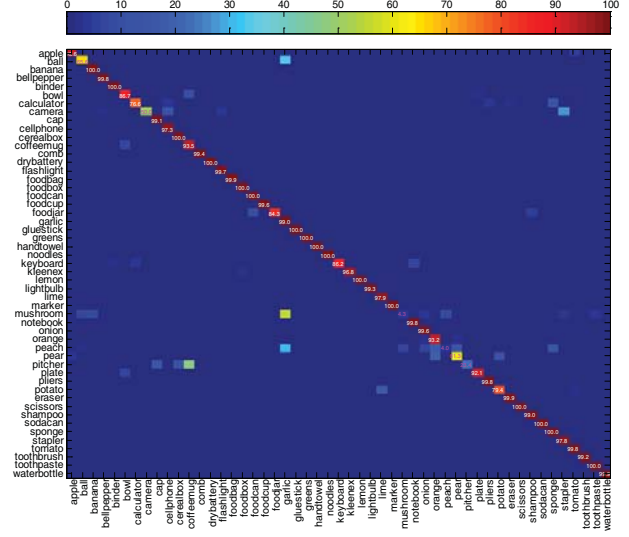


Figure 6. Confusion matrix of CFK on the Washington RGB-D dataset (best viewed in the magnified color image). The y-axis indicates the ground true labels of the testing objects, and the x-axis indicates the predicted labels. Most categories can be recognized correctly, except *pitchers*, *balls*, *mushrooms* and *peaches*, examples of which are shown in Fig. 7.



Figure 7. Examples from the often misclassified categories by CFK. The *pitcher* is classified as the *coffee mug* due to their similar shape, and the *ball*, the *mushroom* and the *peach* can be misclassified as one instance of the category *garlic* because they nearly own the same colors and shapes.

each instance and reduce the size of the database to 2,808 RGB-D images. Then the database is randomly split into training and testing set, where 82 object instances with a total of 1476 views are for training, and the remaining 74 object instances with 1,332 views are for testing. Note that each object instance can only appear in the training set or the testing set. We also repeat the evaluation for 10 times and report the average results.

Results. The comparison of the results on the 2D3D dataset is shown in Table 2. On this dataset, the improvements of our approach are more remarkable. Towards the depth modality-based object recognition, CFK exceeds the previous state-of-the-art SP-HMP with 5.0% (from 87.6% accuracy to 92.6% accuracy). Meanwhile, CFK outperforms SP-HMP by 6.2% (from 86.3% to 92.5%) for the RGB modality-based object recognition. When combining the both modality, CFK also achieves the best result with 94.6% accuracy.

Methods	Depth	RGB	Combine
MLP [8]	74.6	66.6	82.8
SP-HMP [7]	87.6	86.3	91.0
CFK	92.6 ± 1.6	92.5 ± 1.3	94.6 ± 2.0

Table 2. Comparison of results on the 2D3D object database.

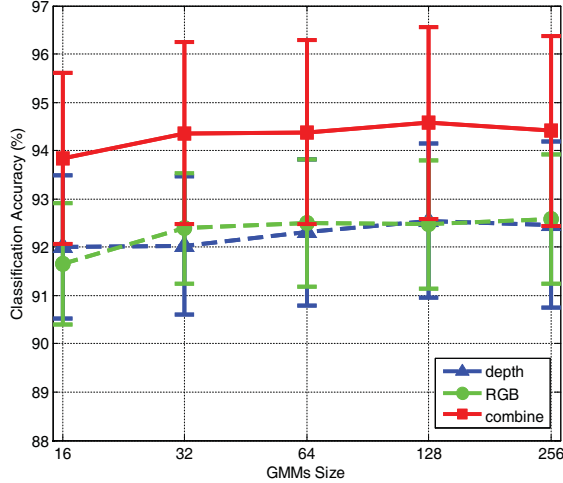


Figure 8. Effect of the GMMs size (M) to the performance of CFK on the 2D3D dataset (2D and 3D spatial pyramids are applied here).

Parameter Analysis. Similarly, we analyze the effects of the GMMs size M and the spatial pyramid matching individually to the performance of CFK on the 2D3D dataset.

(1) The effect of M . Overall, CFK can keep a high recognition accuracy stably when M increases from 16 to 256, as shown in Fig. 8. More detailedly, the larger GMMs size can achieve a few performance gains on this dataset as well, e.g., when GMMs size increases from $M = 16$ to $M = 32$, the corresponding recognition accuracy can be promoted from 93.84% to 94.36%. It is worth noting that even a small size of GMMs, e.g. $M = 16$, are sufficient for CFK to win the best result, compared to all previous methods. This paper fixes $M = 128$ of GMMs on this dataset.

(2) The effect of the 2D/3D spatial pyramid matching. Fig. 9 demonstrates that the 2D SPM (for RGB modality) and the 3D SPM (for depth modality) can clearly improve the recognition performance of CFK, respectively. We conclude that CFK can benefit a lot from the rough spatial structure information to represent an object, including the appearance information of the RGB modality and the shape cue of the depth modality, which agrees with the conclusion of our first experiment as well as the work [25].

Error Analysis. The confusion matrix of CFK over the 14 testing categories of 2D3D dataset is illustrated in Fig. 10. CFK can recognize almost all the categories accurately, but fail to distinguish a few extremely tough examples, as shown in Fig. 11. These misclassified examples are

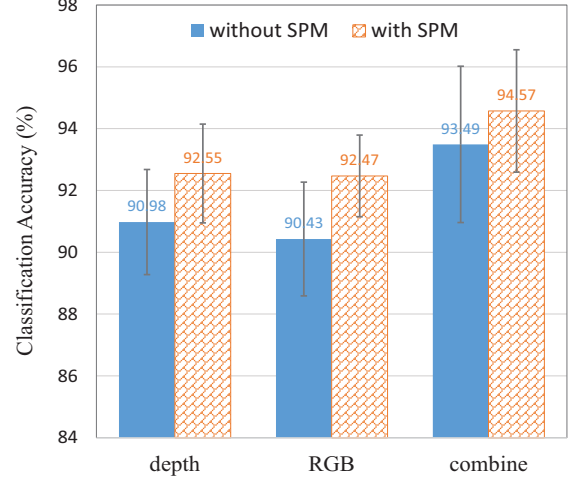


Figure 9. Effect of spatial pyramid matching to the performance of CFK on the 2D3D dataset ($M = 128$).

hard to distinguish because of their easily confused appearance, shape or both with the training examples from other categories. Take the cans for instance, humans are hard to distinguish its appearance from the drink cartoon, if the text labels are ignored.

5. Conclusion and Future Work

This paper proposes a simple yet powerful feature learning method for the novel RGB-D object recognition. We empirically demonstrate that the Fisher Kernel encoding equipped with single-layer convolutional neural networks are sufficient for our CFK approach to achieve the state-of-the-art results on the RGB-D object benchmarks.

Instead of single-layer networks, it is straightforward to stack multiple-layer convolutional neural networks to learn middle level features for each modality of each object (e.g., analogous to the construction of the first layer in Section 3.2, we can build the second layer as well as other layers through learning new filters by sampling the prior convolutional responses. Indeed, such a process has been employed in [7], which proved the hierarchical learning structure were more effective for object representation), and then aggregate them to image-level feature representation by Fisher Kernel encoding with the introduced 2D/3D spatial pyramids. It is worth noting that the proposed CFK approach in this paper belongs to the unsupervised feature learning pipeline. To make use of the labeled data to further improve the performance of CFK, an alternative way is to learn the parameters of CFK in a supervised way. We leave it for future work.

Acknowledgment

This work is funded by the National Basic Research Program of China (Grant No. 2012CB316302), National Nat-

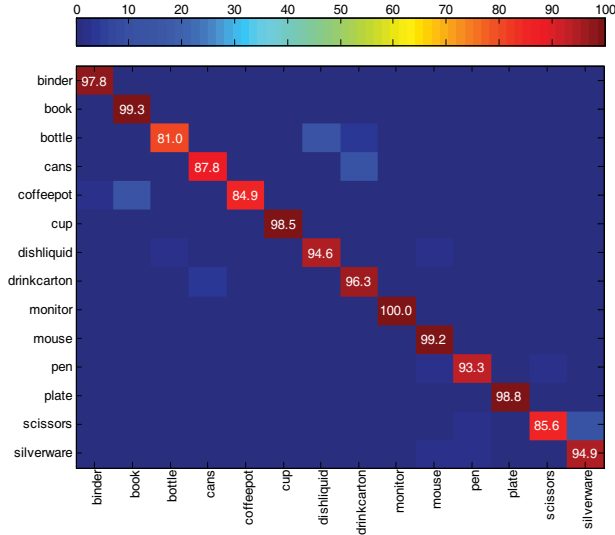


Figure 10. Confusion matrix of CFK on the 2D3D dataset (best viewed in color). Almost all the categories can be recognized correctly, except some very tough examples shown in Fig. 11.

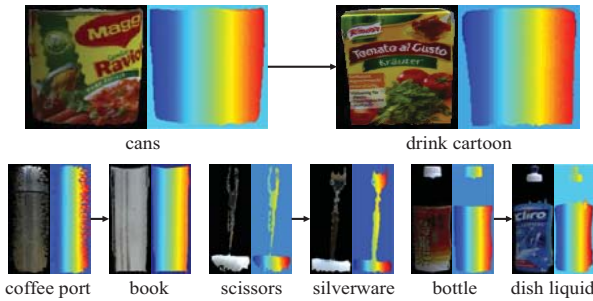


Figure 11. Very tough examples that can be misclassified by CFK. The *cans* has almost the same appearance with the *drink carton*, while the *coffee port*, the *scissors* and the *bottle* are misclassified as the *book*, the *silverware* and the *dish liquid* mainly due to their similar shapes.

ural Science Foundation of China (Grant No. 61322209 and Grant No. 61175007), the Strategic Priority Research Program of the Chinese Academy of Sciences (Grant X-DA06040102).

References

- [1] <https://www.google.com/atap/project-tango/>.
- [2] <https://www.microsoft.com/microsoft-hololens/en-us/>.
- [3] A. Bar-Hillel, D. Hanukaev, and D. Levi. Fusing visual and range imaging for object class recognition. In *ICCV*, 2011.
- [4] M. Blum, J. T. Springenberg, J. Wulfin, and M. Riedmiller. A learned feature descriptor for object recognition in rgb-d data. In *ICRA*, 2012.
- [5] L. Bo, X. Ren, and D. Fox. Depth kernel descriptors for object recognition. In *IROS*, 2011.
- [6] L. Bo, X. Ren, and D. Fox. Hierarchical matching pursuit for image classification: architecture and fast algorithms. In *NIPS*, 2011.
- [7] L. Bo, X. Ren, and D. Fox. Unsupervised feature learning for rgb-d based object recognition. *ISER*, June, 2012.
- [8] B. Browatzki, J. Fischer, B. Graf, H. Bulthoff, and C. Wallraven. Going into depth: Evaluating 2d and 3d cues for object classification on a new, large-scale object dataset. In *ICCV Workshops*, 2011.
- [9] V. L. A. V. Chatfield, Ken and A. Zisserman. The devil is in the details: an evaluation of recent feature encoding methods. In *BMVC*, 2011.
- [10] Y. Cheng, X. Zhao, K. Huang, and T. Tan. Semi-supervised learning for rgb-d object recognition. In *ICPR*, 2014.
- [11] Y. Cheng, X. Zhao, K. Huang, and T. Tan. Semi-supervised learning and feature evaluation for rgb-d object recognition. *Computer Vision and Image Understanding*, 2015.
- [12] A. Coates, A. Y. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. In *AISTATS*, 2011.
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and F.-F. Li. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [14] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge 2012 (voc 2012) results, 2012.
- [15] S. Gupta, R. Girshick, P. Arbelaez, and J. Malik. Learning rich features from rgb-d images for object detection and segmentation. In *ECCV*, 2014.
- [16] Z. W. L. W. Huang, Yongzhen and T. Tan. Feature coding in image classification: A comprehensive study. *PAMI*, 36(3):493–506, 2014.
- [17] A. E. Johnson and M. Hebert. Using spin images for efficient object recognition in cluttered 3d scenes. *PAMI*, 21(5):433–449, 1999.
- [18] I. S. Krizhevsky, Alex and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [19] K. Lai, L. Bo, X. Ren, and D. Fox. A large-scale hierarchical multi-view rgb-d object dataset. In *ICRA*, 2011.
- [20] K. Lai, L. Bo, X. Ren, and D. Fox. Sparse distance learning for object recognition combining rgb and depth information. In *ICRA*, 2011.

- [21] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*, 2006.
- [22] T. Leung and J. Malik. Representing and recognizing the visual appearance of materials using three-dimensional textons. *IJCV*, 43(1):29–44, 2001.
- [23] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, 2004.
- [24] F. Perronnin and C. Dance. Fisher kernels on visual vocabularies for image categorization. In *CVPR*, 2007.
- [25] J. S. Perronnin, Florent and T. Mensink. Improving the fisher kernel for large-scale image classification. In *ECCV*, 2010.
- [26] C. Redondo-Cabrera, R. J. Lopez-Sastre, J. Acevedo-Rodriguez, and S. Maldonado-Bascon. Surfing the point clouds: selective 3d spatial pyramids for category-level object recognition. In *CVPR*, 2012.
- [27] C. Redondo-Cabrera, R. J. Lopez-Sastre, J. Acevedo-Rodriguez, and S. Maldonado-Bascon. Recognizing in the depth: selective 3d spatial pyramid matching kernel for object and scene categorization. *Image and Vision Computing*, 2014.
- [28] S. Shalev-Shwartz, Y. Singer, N. Srebro, and A. Cotter. Pegasos: Primal estimated sub-gradient solver for svm. In *ICML*, 2007.
- [29] A. V. Simonyan, Karen and A. Zisserman. Deep fisher networks for large-scale image classification. In *NIPS*, 2013.
- [30] R. Socher, B. Huval, B. Bath, C. D. Manning, and A. Ng. Convolutional-recursive deep learning for 3d object classification. In *NIPS*, 2012.
- [31] M. S. Sydorov, Vladyslav and C. H. Lampert. Deep fisher kernels—end to end learning of the fisher kernel gmm parameters. In *CVPR*, 2014.
- [32] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong. Locality-constrained linear coding for image classification. In *CVPR*, 2010.
- [33] X. Zhao, Y. Yu, Y. Huang, K. Huang, and T. Tan. Feature coding via vector difference for image classification. In *ICIP*, 2012.
- [34] X. Zhou, K. Yu, T. Zhang, and T. S. Huang. Image classification using super-vector coding of local image descriptors. In *ECCV*, 2010.