

Large-Margin Multi-Modal Deep Learning for RGB-D Object Recognition

Anran Wang, *Student Member, IEEE*, Jiwen Lu, *Senior Member, IEEE*, Jianfei Cai, *Senior Member, IEEE*, Tat-Jen Cham, and Gang Wang, *Member, IEEE*

Abstract—Most existing feature learning-based methods for RGB-D object recognition either combine RGB and depth data in an undifferentiated manner from the outset, or learn features from color and depth separately, which do not adequately exploit different characteristics of the two modalities or utilize the shared relationship between the modalities. In this paper, we propose a general CNN-based multi-modal learning framework for RGB-D object recognition. We first construct deep CNN layers for color and depth separately, which are then connected with a carefully designed multi-modal layer. This layer is designed to not only discover the most discriminative features for each modality, but is also able to harness the complementary relationship between the two modalities. The results of the multi-modal layer are back-propagated to update parameters of the CNN layers, and the multi-modal feature learning and the back-propagation are iteratively performed until convergence. Experimental results on two widely used RGB-D object datasets show that our method for general multi-modal learning achieves comparable performance to state-of-the-art methods specifically designed for RGB-D data.

Index Terms—Deep learning, large-margin feature learning, multi-modality, RGB-D object recognition.

I. INTRODUCTION

RECOGNIZING commonplace objects is a challenging task as real-world scenes are often cluttered and have variable illumination. With the recent advent of commodity depth cameras, an increasing amount of visual data not only contains color but also depth measurements. It is expected that the additional depth information will lead to improved

object recognition performance due to the robustness of depth measurements to light and color variation.

The various approaches that have been proposed for RGB-D object recognition can be divided into two categories from the perspective of feature generation: methods with hand-crafted features [5], [9], [30], and methods with learned features [4], [6], [7], [42], [48], [50]. The hand-crafted or learned features are typically fed into classifiers such as linear SVMs and random forests [8] as done in [4], [30] for the final classification.

In the first category, hand-crafted features such as SIFT [35], SURF [3], textron [34] and color histogram [1] are used to describe color and texture information of color images and 3D geometry information of depth images. The problem with hand-crafted features is that they do not readily extend to different datasets or other modalities, since they are often manually tuned for the conditions encountered in the datasets in mind. In addition, hand-crafted features can only capture a subset of the cues that are useful for recognition.

In the second category, features are learned from raw data for RGB-D object recognition. Representative methods include convolutional-recursive deep learning [42], hierarchical matching pursuit [6], [7], convolutional k-means descriptors [4], hierarchical sparse coding [50] and local coordinate coding [48]. Most of these methods either learn the features for individual modalities independently, or treat RGB-D simply as undifferentiated four-channel data, which cannot adequately exploit the complementary relationship between the two modalities.

To address the above issues, in this paper we propose a general CNN based multi-modal learning method for RGB-D object recognition. The basic idea of our proposed approach is illustrated in Fig. 1. In particular, we build deep CNNs to learn feature representations for color and depth separately, which are then connected with our proposed final multi-modal layer. This layer is designed to not only discover the most discriminative features for each modality, but also harness the complementary relationship between the two modalities. Thus if there is noise in one modality that is uncorrelated with the other modality, the learned features will not include those dimensions particularly susceptible to noise. The results of the multi-modal layer are back-propagated to update the parameters of CNNs and the multi-modal feature learning and the back-propagation are iteratively performed until convergence.

The major contribution of this paper lies in the proposed multi-modal deep learning framework which exploits the complementary information between different modalities, instead

Manuscript received March 13, 2015; revised July 06, 2015 and August 05, 2015; accepted August 07, 2015. Date of publication September 11, 2015; date of current version October 20, 2015. This work was supported by the Singapore National Research Foundation under its International Research Centre@Singapore Funding Initiative and administered by the IDM Programme Office. This work was supported in part by the Singapore Ministry of Education (MOE) Tier 1 RG 138/14, MOE Tier 2 ARC28/14, and the Singapore A*STAR Science and Engineering Research Council under Grant PSF1321202099. The guest editor coordinating the review of this manuscript and approving it for publication was Prof. Benoit Huet.

A. Wang, J. Cai, and T.-J. Cham are with the School of Computer Engineering, Nanyang Technological University (NTU), Singapore 639798 (e-mail: wang001@e.ntu.edu.sg; asjfc@ntu.edu.sg; astjcham@ntu.edu.sg).

J. Lu is with the Department of Automation, Tsinghua University, Beijing 100084, China (e-mail: elujiwen@gmail.com).

G. Wang is with the School of Electrical and Electronic Engineering, Nanyang Technological University (NTU), Singapore 639798 (e-mail: wanggang@ntu.edu.sg).

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

Digital Object Identifier 10.1109/TMM.2015.2476655

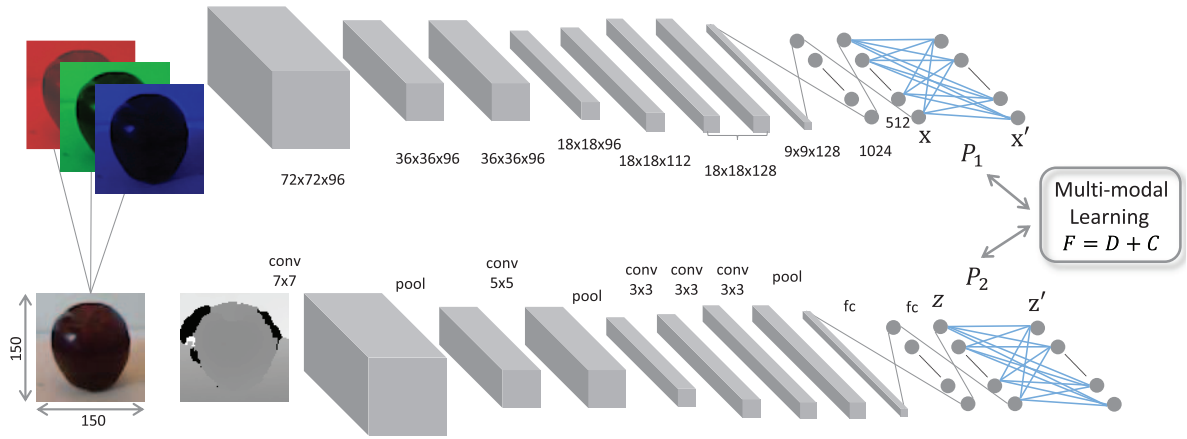


Fig. 1. Overview of our proposed general CNN-based multi-modal learning framework. We first build deep CNNs to learn feature representations of color and depth, and then connect them with a carefully designed multi-modal layer, where two transformation matrices, P_1 and P_2 , are learned to project the activations, x and z , to a new feature space. The cost function of the multi-modal learning contains the discriminative constraints D of each modality and the correlation constraint C between two modalities. Here, “conv”, “pool”, and “fc” stands for “convolutional”, “max-pooling”, and “fully-connected layers”, respectively.

of just treating them as multi-channel input data or concatenating features independently learned from them. The proposed method is a general framework which could be used for other multi-modal applications. In addition, we conduct experiments on two widely used RGB-D object datasets: RGB-D Object Dataset [30] and the 2D3D Dataset [9]. The experimental results show that our general multi-modal learning method with much lower dimensions of features achieves comparable performance to state-of-the-art methods that are specifically designed for RGB-D data.

The rest of this paper is organized as follows. Section II introduces related works on RGB-D object recognition and deep learning. Section III describes our multi-modal deep learning method in detail. Section IV gives the experimental results with discussions and Section V concludes the paper.

II. RELATED WORK

A. RGB-D Object Recognition

Hand-crafted features-based: Past research in image-based object recognition mainly focused on hand designed features such as SIFT and SURF for feature representation. The availability of RGB-D data enabled richer representations to be extracted. For example, Lai *et al.* [30] used hand-crafted features including spin images [25] and SIFT descriptors for depth images, while textons, color histograms, and SIFT descriptors were used for color images. They utilized a bag-of-words model based on these local features and computed an efficient match kernel for object matching. With these encoded features, they evaluated the recognition performance with different classifiers: linear support vector machine, gaussian kernel support vector machine and random forest. Bo *et al.* [5] proposed to combine multiple features such as 3D shape, size features and depth edge features to further improve the recognition performance.

Hand-crafted features have also been extracted by several methods with RGB-D data, but for a different application: RGB-D scene labeling. Silberman *et al.* [40] proposed to use 2D and 3D location prior features besides SIFT and spin image. Koppula *et al.* [28] extracted features such as HSV

color values, HOG, planariness, scatter, vertical component of the normal, vertical position of centroid etc. from superpixels and captured edge features for neighboring superpixels such as angle between normals, horizontal distance between centroids. They achieved excellent performance on scene labeling of single RGB-D images and point clouds reconstructed from RGB-D images. With only depth data, Song *et al.* [43] proposed a robust action recognition framework with a new feature: the body surface context.

Feature learning-based: Recently, several methods have been proposed to learn features from raw data directly in an unsupervised manner for RGB-D object recognition. For example, Bo *et al.* [7] proposed to learn dictionaries via K-SVD [2] from multiple channels, including not only RGB and depth but also gray-scale intensities and surface normals. Then a Hierarchical Matching Pursuit (HMP) method [6] was introduced to generate higher level representations from learned sparse codes of local patches. By using inputs with more channels and two-layer sparse codes, this method captured hierarchical features. Blum *et al.* [4] proposed a feature learning approach from RGB-D data by extracting interest points with SURF features and conducting k-means clustering on image patches around interest points to learn a codebook. Socher *et al.* [42] proposed to use the output of a single-layer CNN that is pre-trained in an unsupervised manner as the input to Recursive Neural Networks (RNNs), so that the CNN produces low-level features while the RNNs learn higher-level features. Features from color and depth channels are learned separately and then concatenated for the final soft-max classifier.

For existing methods that learn features from color and depth modalities separately [42] and concatenate them prior to classification, the major shortcoming is that the relation between the two modalities is ignored and the feature learning stage of one modality would not be effected by other modalities, and thus the complementary nature of different modalities cannot be fully exploited. For other methods that simply combine the modalities from the outset [4], [6], [7] and adopt a conventional approach to learn features, the major shortcoming is that the combination is not physically meaningful and does not capitalize

on the different characteristics of the modalities. Recently, Jhuo *et al.* [23] proposed a multi-modal unsupervised feature learning for RGB-D image classification based on [31]. Their method built a deep structure to extract features from local patches of both gray-scale and depth images followed by spatial pooling and local contrast normalization.

Several methods have been proposed for objects classification with multi-modal information, but with different sources. For example, Izadinia *et al.* [22] proposed a method to segment and localize objects from audio-visual videos. Lu *et al.* [36] and Wu *et al.* [49] presented to classify web multimedia objects with both images and textual descriptions to improve the traditional bag-of-visual-words. Hu *et al.* [21] proposed to learn distance metrics for face verification from multiple features. There are also some works [10], [17], [45] on extracting discriminative information for the 3D model or image retrieval task.

B. Deep Learning

A number of deep learning methods have been proposed in computer vision in recent years, which aim to learn invariant and hierarchical feature representations for different applications. Representative deep learning methods include deep belief networks [19], convolutional deep belief networks [33], deep Boltzmann machines [38], and stacked denoising autoencoders [47]. Most of them have demonstrated promising performance in visual analysis.

Among these deep learning methods, deep convolutional neural networks have produced excellent performance in a number of vision tasks. In 1998, LeCun *et al.* [32] proposed a convolutional neural network structure to deal with handwritten digit recognition problem. In 2012, another representative work was proposed by Krizhevsky *et al.* in [29]. They achieved superior image classification performance on the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) [12], [13] with a large CNN using 1.2 million labeled images. Techniques such as dropout, data augmentation and engineering trick of adjusting step size help to make CNN more trainable. Since then, tremendous interests have been attracted to this deep learning field. CNNs provide consistently high performance on various vision tasks such as scene labeling [15], object detection [16], video classification [27], house number digit classification [39], image super-resolution [14], image quality assessment [26], action recognition [24], face verification [44], pedestrian detection [52], and pose estimation [46]. In addition, deep convolutional neural networks start to show the advancement on text understanding from word-level [37] or character-level [53] inputs. There are also a few methods [51], [41] which visualize the intermediate features of CNNs to help us understand the insight of what CNNs learn and adjust the network design accordingly.

Recently, deep learning has also been used for object recognition and scene labeling in RGB-D data. For example, Couprie *et al.* [11] presented a multi-scale CNN framework for RGB-D scene labeling. They treated RGB-D data as four-channel input to CNNs. Gupta *et al.* [18] proposed a CNN based approach which replaces the original depth map with three channels (horizontal disparity, height above ground, angle between point

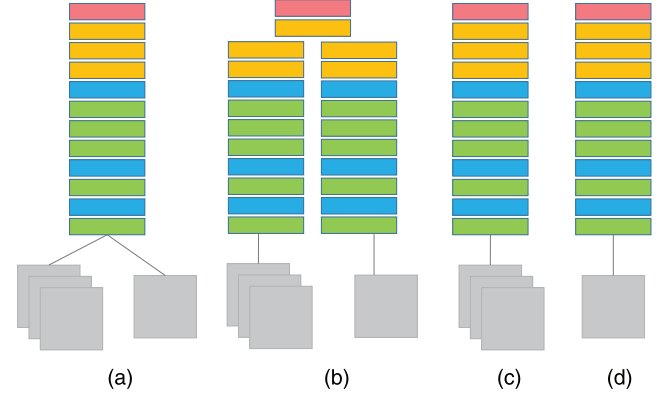


Fig. 2. Different CNN structures for RGB-D data. Green, blue, yellow, and red boxes indicate convolutional, pooling, fully-connected, and softmax layers, respectively.

normal and inferred gravity) as the CNN input for RGB-D object detection and segmentation. They concatenated the features independently learned from color and depth images before the final classification stage. In these deep models, the relationship between color and depth was not adequately exploited.

III. PROPOSED APPROACH

In this section, we describe our proposed CNN-based learning structure followed by detailed formulation, optimization and other implementation specifics.

A. CNN-Based Learning Structure

Of the various ways of using CNNs with RGB-D data, a straightforward approach is to combine RGB and depth data from the outset as a four-channel input to the convolutional neural network, as shown in Fig. 2(a), akin to what is used in Couprie *et al.* [11] for scene labeling. Green, blue, yellow and red boxes indicate convolutional, pooling, fully-connected and softmax layers respectively. Alternatively, discriminative features may be extracted independently from color and depth images by concatenating the activations of the second fully-connected layers of the two modalities, and feeding them into the last fully-connected layer with dense connections. From the final softmax layer, supervised information is back-propagated to the independent networks for both modalities. Such a structure is shown in Fig. 2(b).

Rather than adopting above two approaches that involve a simplistic fusion of color and depth data, in this paper we propose to further explore the relationship between two modalities. We investigate an architecture for multi-modal feature learning carried out in conjunction with convolutional neural networks. In particular, we first pre-train CNNs on color and depth images separately, as shown in Fig. 2(c) and (d). Next the activations of the second fully-connected layers of the two modalities are used in our proposed multi-modal feature learning mechanism, in which not only the most discriminative features are recovered for each modality, but also we exploit the complementary relationship between the two modalities, for which we expect the semantic data to be strongly related but noise sources to be mostly independent. This is achieved by introducing a criterion, related to canonical correlation analysis (CCA) [20], such that

the learned representations are robust to noise that is uncorrelated between the two modalities.

The results of multi-modal learning will then be back-propagated to the lower layers of CNN. The multi-modal feature learning and the back-propagation are iteratively performed until convergence.

B. Notation

Let $X = [x_1, x_2, \dots, x_N] \in \mathbb{R}^{d_1 \times N}$ denote the activations (with d_1 dimensions) of the second fully-connected layer of CNN from color images in one data batch, where there are N images. Similarly, let $Z \in \mathbb{R}^{d_2 \times N}$ be the activations of the same CNN layer from depth images in one data batch. Denote $P_k \in \mathbb{R}^{d_k' \times d_k}$ the transformation matrix for modality k , $k = 1, 2$. Our objective is to learn new representations $x'_i = P_1 x_i$ and $z'_i = P_2 z_i$ which maximize the correlations of color and depth features as well as forcing the distance between same-class objects to be small and the distance between different-class objects to be large. The weights w_1 and w_2 for x'_i and z'_i respectively are also simultaneously learned. The output of our framework will be weighted combination of x'_i and z'_i . Linear SVMs will then be trained to classify these combined features.

C. Formulation

To learn P_k for the two modalities, we consider both discriminative constraints of each modality and the correlation between them. We formulate the cost function as

$$\begin{aligned} \min_{\{P_1, P_2, w_1, w_2\}} \quad & F = w_1 D_1(P_1) + w_2 D_2(P_2) + \lambda C(P_1, P_2) \\ \text{subject to} \quad & w_1 + w_2 = 1, w_1 \geq 0, w_2 \geq 0, \lambda > 0 \end{aligned} \quad (1)$$

where D_k is the discriminative term for modality k , C is the correlation term, and λ is the weight between the discriminative and the correlated constraints.

Discriminative Term: Here we elaborate on the discriminative term for color modality, and likewise for depth modality. The intention is for P_1 to map X to a space in which the distance between x_i and x_j is small if they are of the same class, and large if they are in different classes. Thus we define a constraint as: if two objects are from the same class, their relative feature distance should be smaller than a given threshold $\mu_1 - \tau_1$ ($\mu_1 > \tau_1 > 0$); otherwise, the distance should be larger than $\mu_1 + \tau_1$. Mathematically, this is expressed as

$$y_{ij}(\mu_1 - d_{P_1}(x_i, x_j)) > \tau_1 \quad (2)$$

where labels $y_{ij} = 1$ and $y_{ij} = -1$ indicate x_i and x_j are from the same class and different classes respectively; $d_{P_1}(x_i, x_j)$ is the computed distance; τ_1 and μ_1 are two parameters set empirically. The distance between a pair of the CNN activations, x_i and x_j , in color modality is computed as

$$d_{P_1}(x_i, x_j) = (P_1 x_i - P_1 x_j)^T \cdot (P_1 x_i - P_1 x_j). \quad (3)$$

This leads to the discriminative term defined as

$$D_1(P_1) = \sum_{ij} h(\tau_1 - y_{ij}(\mu_1 - d_{P_1}(x_i, x_j))) \quad (4)$$

where h is a hinge loss function $h(x) = \max(0, x)$. Likewise for P_2 involving depth features.

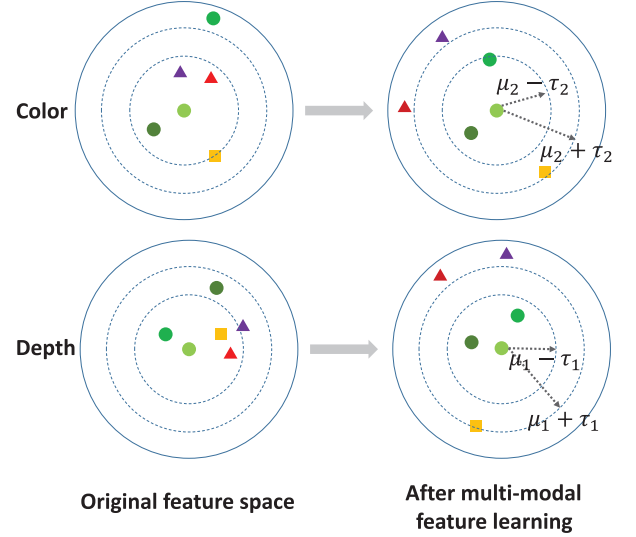


Fig. 3. Illustration of the multi-modal feature learning. Left: objects in the original feature space (CNNs activations). Right: the corresponding objects in the new feature space after the multi-modal learning. Objects dots with the same shape are from the same object class. With the constraints: 1) the distance between different-class objects should be larger than $\mu + \tau$; 2) the distance between same-class objects should be less than $\mu - \tau$; and 3) the difference between the corresponding distances in different modalities should be minimized, we aim at getting more discriminative and complementary features of the two modalities.

Correlation Term: With the discriminative term, features of samples in different classes are better separated for each modality. However, as the data captured in different modalities may suffer from missing information or noise pollution, we need to exploit the relationship between different modalities to reduce misclassification. Following Canonical Correlation Analysis (CCA), we aim to enforce the correlation between the two modalities. One way to maximize the correlation is minimizing the following cost function:

$$C'(P_1, P_2) = \sum_{ij} (P_1 x_i - P_2 z_i)^T \cdot (P_1 x_i - P_2 z_i). \quad (5)$$

For $i = 1, 2, \dots, N$, if $P_1 x_i = P_2 z_i$, the projected features of two modalities are perfectly correlated. But the projected features from different modalities could be of different dimensions, and the difference between them could not be computed directly. Thus we use a relaxation here: enforcing the correlation by minimizing the difference of pairwise distances between the color and depth modalities for all object instances in the dataset, i.e.

$$C(P_1, P_2) = \sum_{ij} \left(\sqrt{d_{P_1}(x_i, x_j)} - \sqrt{d_{P_2}(z_i, z_j)} \right)^2. \quad (6)$$

With such correlation term in addition to the discriminative term, the differences of samples in the feature spaces in different modalities will be constrained to be consistent. In this way, the final features of two modalities are more correlated, and samples of different classes could be separated better with information of two modalities.

Fig. 3 illustrates how these constraints work to get more discriminative and complementary features of the two modalities.

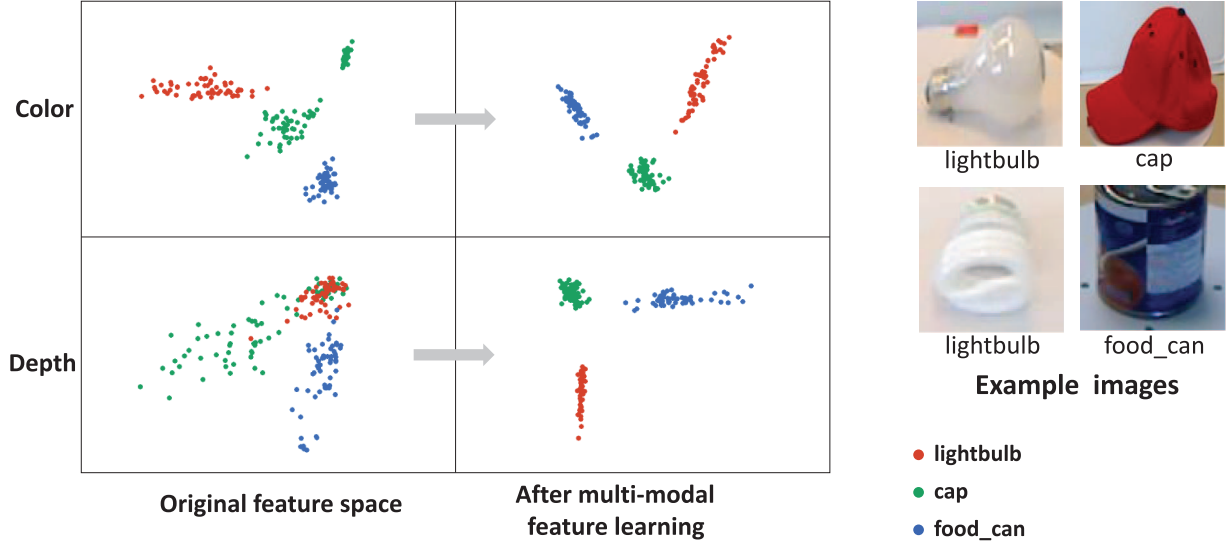


Fig. 4. Examples of the multi-modal feature learning. We show object examples of three classes in both original feature space and learned feature space with PCA projection. Examples of lightbulb, cap, and food_can have similar depth features in the original feature space. With both discriminative and correlation terms, three classes are well separated in the learned feature space.

Object dots with the same shape are from the same object class. The red and purple triangle dots have similar depth features with the yellow square dot in the original feature space. By making sure the pairwise distances of samples in two modalities consistent, the triangle object dots are better separated from the square class in the learned feature space. The correlation term leverages information from two modalities resulting in more complementary features. Fig. 4 shows real examples of objects of RGB-D Object Dataset.

D. Optimization

To our knowledge, there is no closed-form solution to (1) because we need to solve w_k and P_k jointly. To address this, we adopt an alternating approach to optimizing P_k and w_k . In the first stage, the P_k are fixed while w_k are optimized, while in the second stage the opposite occurs. The two stages are alternated until convergence.

When P_k are fixed and w_k is to be updated, we can construct the following Lagrange function based on (1):

$$L(w, \eta) = w_1 D_1 + w_2 D_2 + \lambda C - \eta(w_1 + w_2 - 1). \quad (7)$$

Unfortunately, the solution to (7) will be trivial. When keeping P_1 and P_2 constant, if D_1 is less than D_2 , then the solution to w is: $w_1 = 1$ and $w_2 = 0$, which means the discriminative term for only one modality is effective. It then fails to exploit complementary information of two modalities, and experimentally we found that this leads to suboptimal results. To enable the discriminative terms for both modalities, we make a relaxation of the cost function to

$$\begin{aligned} \min_{\{P_1, P_2, w_1, w_2\}} \quad & F = w_1^p D_1(P_1) + w_2^p D_2(P_2) + \lambda C(P_1, P_2) \\ \text{subject to} \quad & w_1 + w_2 = 1, w_1 \geq 0, w_2 \geq 0, \lambda > 0 \end{aligned} \quad (8)$$

where w_k is replaced by w_k^p . $p > 1$ is an additional parameter. By adding p , the objective becomes nonlinear for w_i and two modalities will be assigned with relatively balanced weights. Thus each modality will have contribution to the final feature representation.

The Lagrange function then becomes

$$L(w, \eta) = w_1^p D_1 + w_2^p D_2 + \lambda C - \eta(w_1 + w_2 - 1). \quad (9)$$

By setting $\frac{\partial L(w, \eta)}{\partial w}$ and $\frac{\partial L(w, \eta)}{\partial \eta}$ to 0, w_k can be updated as

$$w_k = \frac{(1/D_k)^{1/(p-1)}}{\sum_{k=1}^2 (1/D_k)^{1/(p-1)}}. \quad (10)$$

In the second stage we fix w_k and optimize P_k . The optimization is conducted by taking derivative of F with respect to P_k , e.g.

$$\begin{aligned} \frac{\partial F}{\partial P_1} = & 2P_1[w_1^p \sum_{i,j} y_{ij} h'(\tau_1 - y_{ij}(\mu_1 - d_{P_1}(x_i, x_j))) A_{ij}^1 \\ & + \lambda \sum_{i,j} (1 - \sqrt{\frac{d_{P_2}(z_i, z_j)}{d_{P_1}(x_i, x_j)}}) A_{ij}^1] \end{aligned} \quad (11)$$

where $A_{ij}^1 = (x_i - x_j)(x_i - x_j)^T$, and likewise for $\frac{\partial F}{\partial P_2}$. P_k will then be updated using the gradient descent rule

$$P_k = P_k - \beta \frac{\partial F}{\partial P_k}. \quad (12)$$

E. Back-Propagation

In our deep structure, the discriminative information of each modality and the correlation between different modalities are back-propagated to the earlier stage CNNs. Guided by the multi-

feature learning framework, we expect the CNNs to become more informative.

Given the optimized P_k and w_k , the back-propagation is conducted by taking derivative of F with respect to x_i and z_i with

$$\begin{aligned}\frac{\partial D_1}{\partial x_i} &= \sum_j y_{ij} (P_1^T P_1 + P_1 P_1^T) \\ &\quad \times (x_i - x_j) h'(\tau_1 - y_{ij}(\mu_1 - d_{P_1}(x_i, x_j))) \\ \frac{\partial C}{\partial x_i} &= \sum_j (P_1^T P_1 + P_1 P_1^T) \\ &\quad \times (x_i - x_j) \sqrt{\frac{d_{P_1}(x_i, x_j) - d_{P_2}(z_i, z_j)}{d_{P_1}(x_i, x_j)}}. \quad (13)\end{aligned}$$

The same applies for z_i . The whole pipeline of the algorithm is shown in Algorithm 1.

Algorithm 1: Proposed CNN-based multi-modal learning.

Input: Raw RGB and depth images
Output: Transformed CNN features of two modalities: $P_1 x_i, P_2 z_i$, weights between them: w
 //Initialization of CNNs.
 Initialize CNNs by training them with color and depth images respectively as shown in Fig. 2(c) and 2(d)
repeat
 //Alternatively update P_k and w .
 for $k = 1 : 2$ **do**
 Fix P_k :
 Update w according to (10).
 Fix w :
 Update P_k according to (12).
 end for
 //Back-propagation to CNNs.
 Fix P_k and w :
 Update parameters in CNNs according to (13).
until Reaching maximum number of iterations
Return: Optimized P_1, P_2, w

F. Implementation Details

The CNNs are initialized by placing input images in random order and packing them into data batches. Considering all pairwise constraints would require too much computation. Thus, for each iteration, we randomly choose N RGB-D images from all training data and pass them into our multi-modal deep structure. As labels are available, we consider all the pairwise constraints within one batch. After optimizing P_k and w , the CNNs are updated via back-propagation. For simplicity, the convolutional and pooling layers are fixed, and only the fully-connected layers are updated. Each P_k is initialized as an identity matrix, while w is initialized as $[0.5, 0.5]$. For the correlation term, we normalize the distances of different modalities to the same scale. Furthermore in our implementation we use images with two different resolutions as input to capture richer information of different levels. By training with images of two resolutions and

combining their features with equal weight, we find this further improves the performance.

IV. EXPERIMENTS

To evaluate the effectiveness of our proposed multi-modal deep learning method for object recognition, we perform object recognition experiments on the RGB-D Object Dataset [30] and the 2D3D Dataset [9]. The details of the experiments and the results are described in the following sections.

A. Datasets and Experiment Setup

The RGB-D Object Dataset: This dataset has 51 object classes and contains RGB-D images of 300 distinct objects taken from multiple views. These are commonplace objects, some examples of which are cups, keyboards, fruits and vegetables. Each object is video-recorded with cameras mounted at three different heights. There are in total 207,920 RGB-D image frames, with roughly 600 images per object.

Our experiments focused on category recognition. We adopted the same setup as [30] and ran the 10 random splits provided. For each split, one object from each class was sampled, resulting in 51 test objects. After subsampling every 5th frame from the videos, there were some 34,000 images for training and 6900 images for testing.

2D3D Dataset: This dataset consists of 154 objects in 14 different classes. Each object was recorded by a PMD^{TM} Cam-Cube 2.0 time-of-flight camera with views at every 10° around the vertical axis, resulting in a total of 5544 RGB-D images. For category recognition, we adopted the setting of [9]. After excluding some classes with few examples, 6 objects of each class were used for training and the remaining objects were used for testing. For each training (testing) object, only 18 views out of 36 views were used. Eventually 82 objects in 1476 RGB-D images were regarded as training data, while 74 objects in 1332 RGB-D images were used for testing.

Architecture of CNNs: We employed both high resolution and low resolution versions of the RGB-D data, by rescaling all images to 150×150 and 80×80 respectively. For the high-resolution setting of color images, there were 96 kernels of size $7 \times 7 \times 3$ with stride 2, 96 kernels of $5 \times 5 \times 96$ with stride 2, 112 kernels of $3 \times 3 \times 96$ with stride 1, 128 kernels of $3 \times 3 \times 112$ with stride 1, and 128 kernels of $3 \times 3 \times 128$ with stride 1, for the filters of the 1st, 2nd, 3rd, 4th and 5th convolutional layers, respectively. The two fully-connected layers have the sizes of 1024 and 512 respectively. A dropout of 0.5 probability was used for the first fully-connected layer. For each 150×150 images, overlapping 142×142 images were cropped for data augmentation. There were max-pooling layers following the first, the second and the fifth convolutional layers. ReLu non-linearity [29] was applied to the output of every convolutional layer and every fully-connected layer. Note that when initializing CNNs by independently training with color and depth images as shown in Fig. 2(c) and (d), the final fully-connected layer had a size equal to the number of categories, which were then fed into the final softmax layer.

We used the same architecture for both color and depth, apart from the size of filters in the first convolutional layer (3 channels for color and 1 channel for depth). In the low-resolution set-

TABLE I
COMPARISON OF DIFFERENT BASELINES OF
USING CNNs ON RGB-D OBJECT DATASET

Method	Accuracy (%)
RGB CNN	74.6 ± 2.9
Depth CNN	75.5 ± 2.7
RGB-D CNN with 4-channel input	80.2 ± 1.9
RGB-D CNN connected at fc	84.7 ± 2.1

ting, we removed the second convolutional layer and the second pooling layer compared to the high-resolution setting.

Parameters Setting: We set the size of the data batch N to 128. Weights of the CNN convolutional and fully-connected layers were initialized by a zero-mean Gaussian distribution with standard deviation 0.01. The learning rate was first set to 0.01. Once there was no performance improvement on validation data (randomly choosing one object from the testing data for each category), the learning rate was updated to 0.001. For our multi-modal layer, the dimensionality d_k' of the transformed features was set to be the same as d_k , although it could have been different. The parameters $p, \beta, \lambda, \mu_1, \tau_1, \mu_2, \tau_2$ were empirically set as 2, 0.0003, 0.001, 10, 1, 300, 30 respectively for all the experiments for the two resolutions of both the RGB-D Object and 2D3D Datasets.

B. Results on RGB-D Object Dataset

Comparison with different baselines of using CNNs: We conducted experiments on 150×150 input images with four different CNN-based baselines: 1) CNN trained using RGB images only [Fig. 2(c)], named “RGB CNN”; 2) CNN trained using depth images only [Fig. 2(d)], named “Depth CNN”; 3) RGB-D used as the four-channel input to a CNN [Fig. 2(a)], named “RGB-D CNN with 4-channel input”; 4) CNN with separate training for color and depth at the lower layers, followed by concatenating the activations of the second fully-connected layer and feeding them into the last fully-connected layer [Fig. 2(b)], named “RGB-D CNN connected at fc ”.

Table I shows the overall recognition accuracy of the four baselines on the RGB-D object dataset. Observe that although simply adding depth as the fourth channel of the CNN input improved performance, extracting features separately from color and depth and connecting them at the later stage performed significantly better. This is because separately learning features at the early stage for different modalities may result in the features being more independent; otherwise the CNN will primarily learn features for the predominant modality.

Table II compares the results of our proposed multi-model learning with the best baseline, “RGB-D CNN connected at fc ”. We conducted the experiments on the inputs with two resolutions: 150×150 and 80×80 . It can be seen that our method led to improved performance for both cases. This is mainly because our method considers both discriminative constraints of each modality as well as the correlation constraint between two modalities, which will not happen through simply connecting color and depth by a fully-connected layer. By combining the learned features from high-resolution and low-resolution data together, the performance was further improved. In Fig. 5, we show some examples which are correctly classified by our

TABLE II
COMPARISON BETWEEN THE BEST BASELINE METHOD AND OUR METHOD
ON RGB-D OBJECT DATASET WITH INPUTS OF TWO RESOLUTIONS:
 150×150 AND 80×80

Method	Accuracy (%)
RGB-D CNN connected at fc -150	84.7 ± 2.1
Ours-150	85.8 ± 2.6
RGB-D CNN connected at fc -80	81.3 ± 2.4
Ours-80	83.6 ± 2.3
RGB-D CNN connected at fc -150+80	85.6 ± 2.2
Ours-150+80	86.7 ± 2.7



Fig. 5. Examples which are correctly classified by our multi-modal method but misclassified by the baseline with fc connecting two modalities.

multi-modal method but misclassified by the baseline with fc layer connecting two modalities.

Comparison with different settings of training data: Generally, higher resolution images and more data lead to better performance because of more useful information. As we follow the setup as [7], in the above experimental results, we sampled every 5th frame from each video in the dataset. When we use all frames for training as [42], we get an accuracy improvement of 0.2% as shown in Table III. The gain is not so significant due to the high similarity among neighboring video frames as well as the fact that all RGB-D images are sampled from just around 300 videos. We also tried even higher resolution images, but there is no performance gain. The main reason is that the resolutions of original images in the RGB-D dataset are mostly lower or around 150×150 . Besides, even if the resolution of input images can be increased, the network size will be too large to be trained by the limited training data. We believe with more diverse labeled data, we can improve the performance further.

Comparison with state-of-the-art methods: We compared our method with: 1) Lai *et al.* [30]: using SIFT and spin images for depth, and SIFT, color histogram and texton histogram for color; 2) Blum *et al.* [4]: using convolutional k-means descriptors;

TABLE III
COMPARISON WITH RESULT USING MORE LABELED DATA

Method	Accuracy (%)
Ours-150+80	86.7 ± 2.7
Ours-150+80 (more labeled data)	86.9 ± 2.6

TABLE IV
COMPARISON WITH STATE-OF-THE-ART
METHODS ON RGB-D OBJECT DATASET

Method	Accuracy (%)
Lai <i>et al.</i> [30]	81.9 ± 2.8
Blum <i>et al.</i> [4]	86.4 ± 2.3
Socher <i>et al.</i> [42]	86.8 ± 3.3
Bo <i>et al.</i> [7]	87.5 ± 2.9
Le <i>et al.</i> [31]	86.7 ± 2.7
Jhuo <i>et al.</i> [23]	89.6 ± 3.8
Ours	86.9 ± 2.6

3) Socher *et al.* [42]: using Recursive Neural Network + Convolutional Neural Network; 4) Bo *et al.* [7]: using sparse coding based feature learning with additional input channels such as surface normal; 5) Le *et al.* [31]: using features extracted with RICA (Independent Components Analysis with Reconstruction cost) constraints; 6) Jhuo *et al.* [23]: conducting spatial pooling and local contrast normalization after learning features for local patches. Note that the result of [31] on the RGB-D object dataset is reported in [23].

The comparison results are shown in Table IV. It can be seen that our method achieved very competitive results compared with state-of-the-art methods. Note that our method is a general one that can be easily applied to other different modalities beyond color and depth, while the other methods were specifically designed for RGB-D data. Methods of Bo *et al.* [7] and Socher *et al.* [42] have large feature vector sizes: 188300 and 32768 respectively, while ours has only 2048 dimensions, which is much more efficient. Compared with method of Socher *et al.* [42] which also uses Convolutional Neural Network, our method has some fundamental difference: 1) their method has two stages: single-layer CNN feature extraction, random RNNs merging process. Our structure is very compact as it has a uniform structure for different layers; 2) Our method is supervised while [42] is unsupervised. As shown in Table III, although added data from neighboring frames are redundant, there is still some gain in performance. We believe with more diverse labeled data, we can improve the performance further; 3) Their method did not consider inter-modal information. In our method, with both inter-modal and intra-modal constraints, we achieve 1.2% improvement compared with the one only considering the discriminative constraint in each modality.

Compared to the method of Jhuo *et al.* [23] that reports the best classification accuracy, our method has the following two major advantages: 1) [23] has large feature vector sizes: about 394272 for RGB-D Object Dataset and 609408 for 2D3D Dataset. The dimension of our features is 2048, which is much more efficient for computation. 2) For exploiting the relationship between the two modalities, [23] constrains the transformation matrices of different modalities to be identical,

which is a very strict constraint. Our method provides a more general framework which learns one transformation matrix for each modality. In addition, with their constraint, [23] requires the learned features for different modalities to be the same in dimensions while our method could capture more general and diverse information.

The confusion matrix of our final results is shown in Fig. 6, whose diagonal elements represent the recognition accuracy for each category. Fig. 7 shows a few misclassification examples. For instance, the bell pepper was misclassified as a food cup due to their geometrical shapes being very similar. Likewise, the keyboard was misclassified as a kleenex box due to strong similarities in both color and shape.

C. Results on 2D3D Dataset

The number of training samples was too small in the 2D3D Dataset for direct training of CNNs which will lead to severe over-fitting. Thus, we use the CNNs pretrained on the RGB-D Object Dataset for the initialization. Table V shows the performance of different fine-tunings with different CNN architectures shown in Fig. 2. Specifically, two fine-tuning strategies were considered: 1) Fix filters in convolutional layers and only update the parameters of fully-connected layers, named “Update fc ”; 2) Update the parameters of both convolutional and fully-connected layers, named “Update $conv\&fc$ ”. It can be seen from Table V that updating both convolutional and fully-connected layers improved the performance. This is because the data in the two datasets have different characteristics, e.g. images from 2D3D dataset have richer texture. Filters trained on RGB-D Object Dataset might not be suitable for 2D3D Dataset. By adjusting both convolutional and fully-connected layers, the network could better represent the 2D3D dataset.

Table VI shows the comparison between our method and the best baseline approach, “RGB-D CNN connected at fc ”, with inputs of two resolutions. Table VII provides a comparison of our method and the existing methods that reported results on this dataset, including 1) Browatzki *et al.* [9]: using multiple descriptors such as 3D shape context and depth buffer for depth and SURF and self similarity features for color; 2) Bo *et al.* [7]; 3) Le *et al.* [31]; 4) Jhuo *et al.* [23]. Similar remarks as those for the results on the RGB-D Object Dataset can be made here, i.e. our method for general multi-modal feature learning with lower dimensions of features achieves comparable performance to state-of-the-art methods that are specifically designed for RGB-D data or with much higher dimension of features.

D. Other Results

Parameter analysis: With initialized CNNs, we updated our multi-modal learning structure for 300 iterations. In each iteration, 128 randomly chosen samples were passed through the structure to update P_k , w and the parameters of CNN layers. For total 300 iterations, there are around 60,000 positive and 2,370,000 negative pairwise constraints. Fig. 8 shows the value of the objective function under different numbers of iterations on RGB-D Object Dataset with input images of 150×150 . We also show the performance versus different λ in (8) on split 1 of RGB-D Object Dataset in Fig. 9. We can see that an excessively

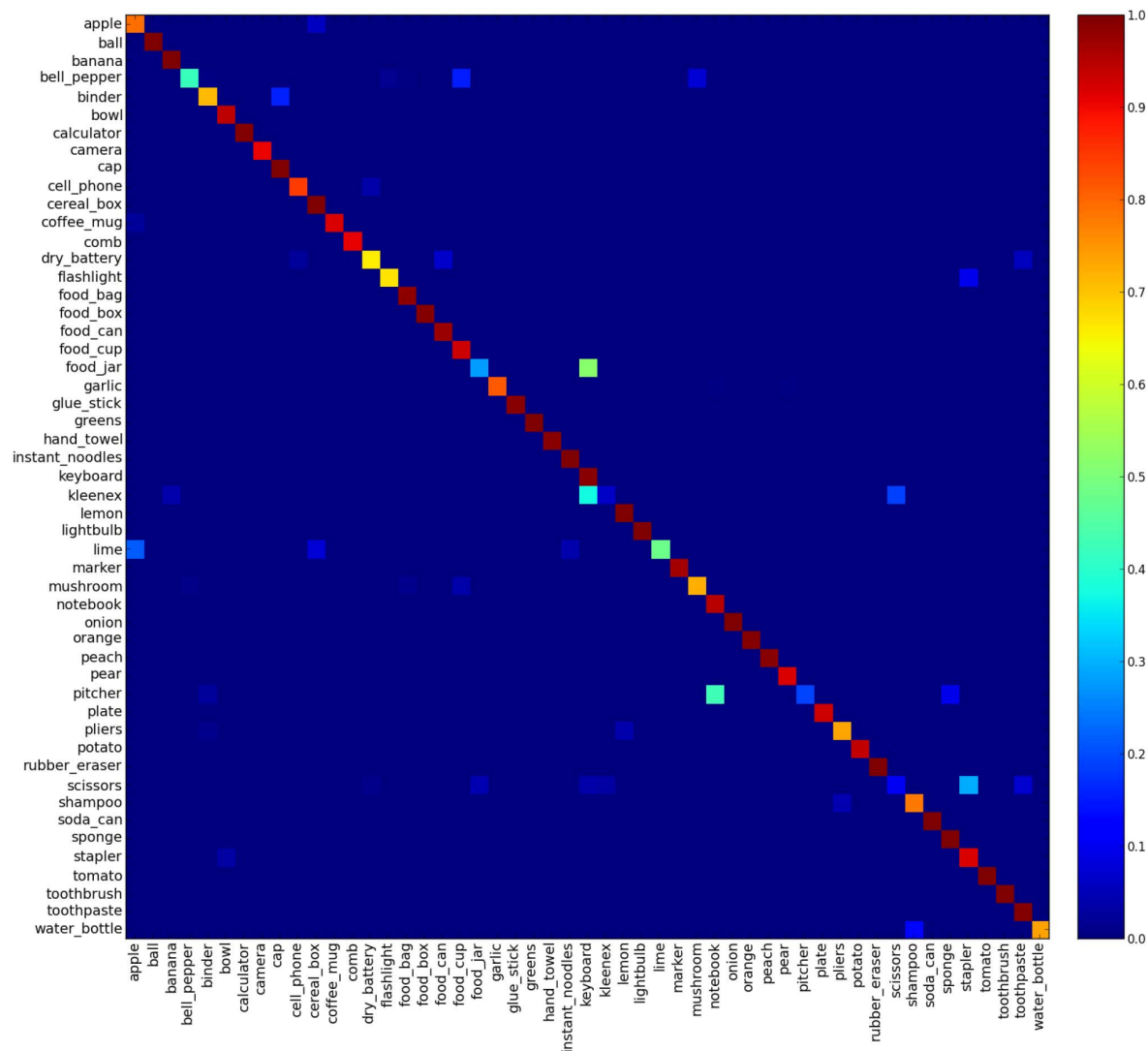


Fig. 6. Confusion matrix of the category recognition results on RGB-D Object Dataset. The vertical axis shows the true labels and the horizontal axis shows the predicted labels.



Fig. 7. Misclassification examples. Misclassifications are due to similar color, texture, or geometry shape.

large or small weight λ for the correlation term led to a drop in performance.

Computational time: We reused the code of [29] on a Titan GPU. For pre-training CNNs with color and depth images separately, stable performance can be reached in roughly 30 epochs.

TABLE V
PERFORMANCE WITH DIFFERENT FINE-TUNING STRATEGIES ON 2D3D DATASET

Accuracy %	RGB CNN	Depth CNN	RGB-D CNN connected at <i>fc</i>
Update <i>fc</i>	56.2	75.4	76.0
Update <i>conv&fc</i>	62.8	85.9	86.5

On the RGB-D Object Dataset with the high resolution setting, pre-training took about 2.5 hours. The optimization of our multi-modal structure required around 2 hours for both datasets. The testing process took less than 1 minute.

Convergence of our algorithm: To show the convergence of our algorithm, we conduct experiment with fixed training data. We fix 100 sampled training batches for each iteration and consider all constraints in each batch. We obtain objective values for 10 iterations, which are decreasing continuously, as shown in Fig. 10.

Different ways of using depth: As shown above, our framework could generally take RGB and depth images as the input

TABLE VI
COMPARISON BETWEEN THE BEST BASELINE METHOD AND
OUR METHOD ON 2D3D DATASET WITH INPUTS OF TWO
RESOLUTIONS: 150×150 AND 80×80

Method	Accuracy (%)
Connect RGB, depth with <i>fc</i> layer-150	86.5
Ours-150	87.4
Connect RGB, depth with <i>fc</i> layer-80	78.4
Ours-80	81.2
Connect RGB, depth with <i>fc</i> layer-150+80	87.1
Ours-150+80	88.4

TABLE VII
COMPARISON WITH OTHER METHODS ON 2D3D DATASET

Method	Accuracy (%)
Browatzki <i>et al.</i> [9]	82.8
Bo <i>et al.</i> [7]	91.0
Le <i>et al.</i> [31]	91.5
Jhuo <i>et al.</i> [23]	92.7
Ours	88.4

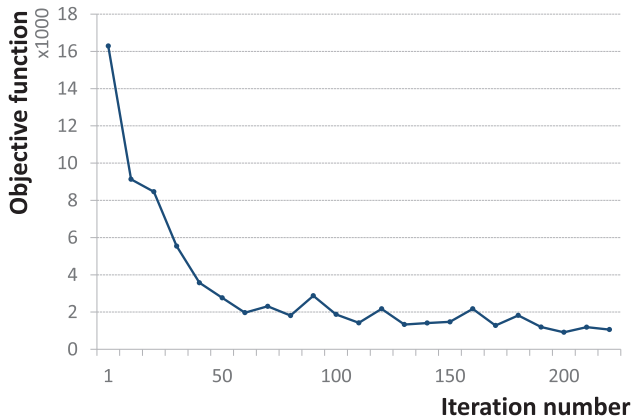


Fig. 8. Value of the cost function under different numbers of iterations.

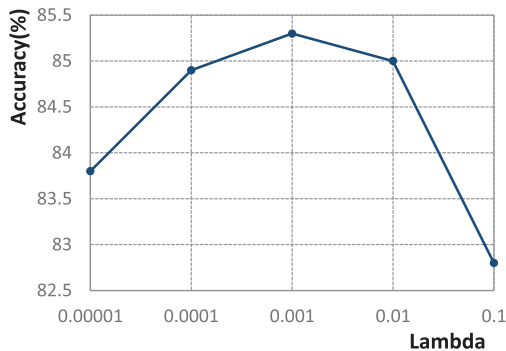


Fig. 9. Effect of choosing different λ .

and provide comparable results on RGB-D object recognition application without specific design. On the other hand, for this particular application, there are different ways of using depth map. As mentioned in the Section II, Gupta *et al.* [18] replaced the original depth map with three channels (horizontal disparity, height above ground, angle between point normal and inferred

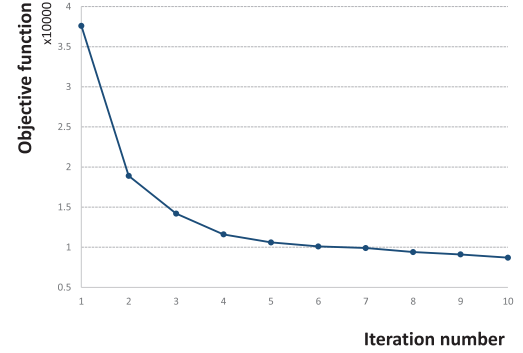


Fig. 10. Convergence of our algorithm.

TABLE VIII
COMPARISON ON USING SURFACE NORMAL INSTEAD OF DEPTH

Method	Accuracy (%)
RGB CNN	74.6 ± 2.9
Depth CNN	75.5 ± 2.7
Surface normal CNN	76.3 ± 2.5
Ours-150 RGB-D	85.8 ± 2.6
Ours-150 RGB+Surface normal	86.2 ± 2.5

gravity) as the CNN input for RGB-D object detection and segmentation. Since in our object recognition application, the information about the ground or gravity is unavailable, we try replacing depth with surface normal. Table VIII shows the results on RGB-D Object Dataset, which indicate that surface normals can better represent geometry information than depth map.

E. Remarks

We make the following key observations from experimental results represented in Tables I–VIII and Fig. 5–10:

- Among different straightforward ways of applying CNNs on RGB-D object recognition, the CNN structure, which separately trains color and depth at the lower layers and concatenates the activations of the higher fully-connected layer, achieved the best performance. It is because independent features could be derived from different modalities using the separate learning at the early stage. Otherwise, features of the predominant modality will be primarily learned.
- Our proposed multi-model learning method outperforms the best baseline of straightforwardly applying CNNs. This is because our method considers not only discriminative constraints of each modality but also the correlation constraint between two modalities.
- We have shown that by using more labeled data, the performance of our method is improved. We believe with more and more diverse training data being added, the performance can be further boosted.
- Different from state-of-the-art methods using hand crafted features, our method directly learns features from raw training data. Compared to state-of-the-art feature-learning based methods, we achieved comparable results with much lower dimension of features. Unlike the existing methods which just concatenate RGB and depth from the outset or concatenate the independently learned features, our

method considers both inter-modal and intra-modal constraints which exploits the complementary multi-modal information better. In addition, our framework is general and can be extended to different multi-modal applications.

- By fine-tuning the pretrained model on a relative large dataset, our method can achieve comparable performance with state-of-the-art methods on dataset with less than 2000 training images.

V. CONCLUSION

In this paper, we proposed a multi-modal deep learning framework to tackle the RGB-D object recognition problem. Compared with the methods that fuse of color and depth at the fully-connected layer, our method simultaneously learns transformation matrices for two modalities with a large margin criterion and a maximal cross-modality correlation criterion. By iteratively optimizing the multi-modal learning framework and updating CNN parameters in conjunction with multi-resolution inputs, we obtained comparable performance to state-of-the-art methods on the RGB-D Object and 2D3D Datasets. While state-of-the-art methods are specifically designed for RGB-D data, our framework is general and can be easily extended to other different modalities.

ACKNOWLEDGMENT

The authors would like to thank J. Hu with the Nanyang Technological University for his helpful discussion.

REFERENCES

- [1] A. E. Abdel-Hakim and A. A. Farag, "CSIFT: A sift descriptor with color invariant characteristics," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2006, vol. 2, pp. 1978–1983.
- [2] M. Aharon, M. Elad, and A. Bruckstein, "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Trans. Signal Process.*, vol. 54, no. 11, pp. 4311–4322, Nov. 2006.
- [3] H. Bay, T. Tuytelaars, and L. Van Gool, "SURF: Speeded up robust features," in *Proc. Eur. Conf. Comput. Vis.*, 2006, pp. 404–417.
- [4] M. Blum, J. T. Springenberg, J. Wulfin, and M. Riedmiller, "A learned feature descriptor for object recognition in RGB-D data," in *Proc. IEEE Int. Conf. Robot. Automat.*, May 2012, pp. 1298–1303.
- [5] L. Bo, X. Ren, and D. Fox, "Depth kernel descriptors for object recognition," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Sep. 2011, pp. 821–826.
- [6] L. Bo, X. Ren, and D. Fox, "Hierarchical matching pursuit for image classification: Architecture and fast algorithms," in *Proc. NIPS*, 2011, pp. 2115–2123.
- [7] L. Bo, X. Ren, and D. Fox, "Unsupervised feature learning for RGB-D based object recognition," in *Proc. ISER*, 2012, pp. 387–402.
- [8] L. Breiman, "Random forests," *Machine Learning*, no. 45, pp. 5–32, 2001.
- [9] 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 *Proc. IEEE Int. Conf. Comput. Vis. Workshop*, Nov. 2011, pp. 1189–1195.
- [10] J.-Y. Chen, C.-H. Lin, P.-C. Hsu, and C.-H. Chen, "Point cloud encoding for 3D building model retrieval," *IEEE Trans. Multimedia*, vol. 16, no. 2, pp. 337–345, Feb. 2014.
- [11] C. Couprie, C. Farabet, L. Najman, and Y. LeCun, "Indoor semantic segmentation using depth information," *CoRR*, 2013 [Online]. Available: <http://arxiv.org/abs/1301.3572>
- [12] O. Russakovsky *et al.*, "ImageNet large scale visual recognition challenge," *Int. J. Comput. Vis.*, pp. 1–42, Apr. 2015 [Online]. Available: <http://www.image-net.org/challenges/LSVRC/2012/>
- [13] J. Deng *et al.*, "ImageNet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2009, pp. 248–255.
- [14] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 184–199.
- [15] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1915–1929, Aug. 2013.
- [16] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2014, pp. 580–587.
- [17] B. Gong, J. Liu, X. Wang, and X. Tang, "Learning semantic signatures for 3D object retrieval," *IEEE Trans. Multimedia*, vol. 15, no. 2, pp. 369–377, Feb. 2013.
- [18] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, "Learning rich features from RGB-D images for object detection and segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 345–360.
- [19] G. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [20] H. Hotelling, "Relations between two sets of variates," *Biometrika*, pp. 321–377, 1936.
- [21] J. Hu, J. Lu, J. Yuan, and Y.-P. Tan, "Large margin multi-metric learning for face and kinship verification in the wild," in *Proc. ACCV*, 2014, pp. 252–267.
- [22] H. Izadinia, I. Saleemi, and M. Shah, "Multimodal analysis for identification and segmentation of moving-sounding objects," *IEEE Trans. Multimedia*, vol. 15, no. 2, pp. 378–390, Feb. 2013.
- [23] I.-H. Jhuo, S. Gao, L. Zhuang, D. Lee, and Y. Ma, "Unsupervised feature learning for RGB-D image classification," in *Computer Vision—ACCV 2014*. New York, NY, USA: Springer, 2015, pp. 276–289.
- [24] S. Ji, W. Xu, M. Yang, and K. Yu, "3D convolutional neural networks for human action recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 1, pp. 221–231, Jan. 2013.
- [25] A. E. Johnson and M. Hebert, "Using spin images for efficient object recognition in cluttered 3D scenes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 21, no. 5, pp. 433–449, May 1999.
- [26] L. Kang, P. Ye, Y. Li, and D. Doermann, "Convolutional neural networks for no-reference image quality assessment," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2014, pp. 1733–1740.
- [27] A. Karpathy *et al.*, "Large-scale video classification with convolutional neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2014, pp. 1725–1732.
- [28] H. S. Koppula, A. Anand, T. Joachims, and A. Saxena, "Semantic labeling of 3D point clouds for indoor scenes," in *Proc. Adv. Neural Inf. Process. Syst.*, 2011, pp. 244–252.
- [29] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [30] K. Lai, L. Bo, X. Ren, and D. Fox, "A large-scale hierarchical multi-view RGB-D object dataset," in *Proc. IEEE Int. Conf. Robot. Automat.*, May 2011, pp. 1817–1824.
- [31] Q. V. Le, A. Karpenko, J. Ngiam, and A. Y. Ng, "Ica with reconstruction cost for efficient overcomplete feature learning," in *Proc. Adv. Neural Inf. Process. Syst.*, 2011, pp. 1017–1025.
- [32] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [33] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," in *Proc. ICML*, 2009, pp. 609–616.
- [34] T. Leung and J. Malik, "Representing and recognizing the visual appearance of materials using three-dimensional textons," *IJCV*, vol. 43, no. 1, pp. 29–44, 2001.
- [35] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *IJCV*, vol. 60, no. 2, pp. 91–110, 2004.
- [36] W. Lu *et al.*, "Web multimedia object classification using cross-domain correlation knowledge," *IEEE Trans. Multimedia*, vol. 15, no. 8, pp. 1920–1929, Dec. 2013.
- [37] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proc. NIPS*, 2013, pp. 3111–3119.
- [38] R. Salakhutdinov and G. E. Hinton, "Deep Boltzmann machines," in *Proc. Int. Conf. Artificial Intell. Statist.*, 2009, pp. 448–455.
- [39] P. Sermanet, S. Chintala, and Y. LeCun, "Convolutional neural networks applied to house numbers digit classification," in *Proc. Int. Conf. Pattern Recog.*, Nov. 2012, pp. 3288–3291.
- [40] N. Silberman and R. Fergus, "Indoor scene segmentation using a structured light sensor," in *Proc. IEEE Int. Conf. Comput. Vis. Workshops*, Nov. 2011, pp. 601–608.
- [41] K. Simonyan, A. Vedaldi, and A. Zisserman, "Deep inside convolutional networks: Visualising image classification models and saliency maps," *CoRR*, 2013 [Online]. Available: <http://arxiv.org/abs/1312.6034>

- [42] R. Socher, B. Huval, B. Bhat, C. D. Manning, and A. Y. Ng, "Convolutional-recursive deep learning for 3D object classification," in *Proc. NIPS*, 2012, pp. 665–673.
- [43] Y. Song, J. Tang, F. Liu, and S. Yan, "Body surface context: A new robust feature for action recognition from depth videos," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 24, no. 6, pp. 952–964, Jun. 2014.
- [44] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, "DeepFace: Closing the gap to human-level performance in face verification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2014, pp. 1701–1708.
- [45] J. Tang, Z. Li, M. Wang, and R. Zhao, "Neighborhood discriminant hashing for large-scale image retrieval," *IEEE Trans. Image Process.*, vol. 24, no. 9, pp. 2827–2840, Sep. 2015.
- [46] A. Toshev and C. Szegedy, "DeepPose: Human pose estimation via deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2014, pp. 1653–1660.
- [47] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and composing robust features with denoising autoencoders," in *Proc. ICML*, 2008, pp. 1096–1103.
- [48] J. Wang *et al.*, "Locality-constrained linear coding for image classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2010, pp. 3360–3367.
- [49] L. Wu, Y. Hu, M. Li, N. Yu, and X.-S. Hua, "Scale-invariant visual language modeling for object categorization," *IEEE Trans. Multimedia*, vol. 11, no. 2, pp. 286–294, Feb. 2009.
- [50] K. Yu, Y. Lin, and J. Lafferty, "Learning image representations from the pixel level via hierarchical sparse coding," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2011, pp. 1713–1720.
- [51] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 818–833.
- [52] X. Zeng, W. Ouyang, M. Wang, and X. Wang, "Deep learning of scene-specific classifier for pedestrian detection," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 472–487.
- [53] X. Zhang and Y. LeCun, "Text understanding from scratch," *CoRR*, 2015 [Online]. Available: <http://arxiv.org/abs/1502.01710>



Anran Wang (S'14) received the B.S. degree from Tianjin University, Tianjin, China in 2012, and is currently working toward the Ph.D. degree in computer engineering at Nanyang Technological University, Singapore.

Her research interests are computer vision and machine learning.



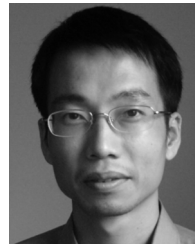
Jiwen Lu (S'10–M'11–SM'15) received the B.Eng. degree in mechanical engineering and the M.Eng. degree in electrical engineering from the Xi'an University of Technology, Xi'an, China, in 2003 and 2006, respectively, and the Ph.D. degree in electrical engineering from the Nanyang Technological University, Singapore, in 2011.

He is an Associate Professor with the Department of Automation, Tsinghua University, Beijing, China. He has authored or co-authored over 110 scientific papers, and more than 40 of those papers were published in IEEE transactions and journals and top-tier computer vision conferences.

His current research interests include computer vision, pattern recognition, and machine learning.

Dr. Lu serves as an Associate Editor of *Pattern Recognition Letters* and *Neurocomputing*. He was a recipient of the First-Prize National Scholarship and the National Outstanding Student Award from the Ministry of Education of China

in 2002 and 2003, the Best Student Paper Award from Pattern Recognition and Machine Intelligence Association of Singapore in 2012, the Top 10% Best Paper Award from the IEEE International Workshop on Multimedia Signal Processing in 2014, and the National 1000 Young Talents Plan Program in 2015.

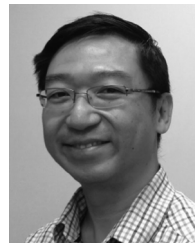


Jianfei Cai (S'97–M'02–SM'07) received the Ph.D. degree from the University of Missouri, Columbia, MO, USA, in 2002.

He is currently an Associate Professor and has served as the Head of Visual & Interactive Computing Division and the Head of Computer Communication Division with the School of Computer Engineering, Nanyang Technological University, Singapore. He has authored more than 150 technical papers in international journals and conferences. His major research interests include

computer vision, visual computing, and multimedia networking.

Dr. Cai has been serving as an Associate Editor for the IEEE TRANSACTIONS ON IMAGE PROCESSING since 2013. He also served as an Associate Editor for the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY from 2006 to 2013. He has been actively participating in program committees of various conferences. He has served as the leading Technical Program Chair for the IEEE International Conference on Multimedia & Expo 2012 and the leading General Chair for the Pacific-Rim Conference on Multimedia 2012.

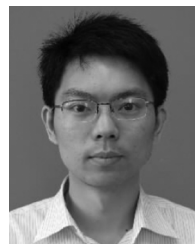


Tat-Jen Cham received the B.A. degree in engineering and the Ph.D. degree from the University of Cambridge, Cambridge, U.K., in 1993 and 1996, respectively.

He received a Jesus College Research Fellowship in Cambridge, U.K., from 1996 to 1997 and was a Research Scientist with the DEC/Compaq Research Lab, Boston, MA, USA, from 1998 to 2001. He was the Director for the Centre of Multimedia and Network Technology, Nanyang Technological University (NTU), Singapore, from 2007 to 2015,

and a Senator with NTU from 2010 to 2014. He is currently an Associate Professor with the School of Computer Engineering, NTU, and a Principal Investigator with the NRF BeingThere Centre for 3D Telepresence, Institute for Media Innovation, NTU. His research interests include computer vision and machine learning.

Prof. Cham was a Singapore–MIT Alliance Fellow from 2003 to 2006. He is on the Editorial Board for the *International Journal of Computer Vision*, and was a General Co-Chair for the Asian Conference on Computer Vision 2014.



Gang Wang (M'11) received the B.S. degree from the Harbin Institute of Technology in Electrical Engineering, Harbin, China, in 2005, and the Ph.D. degree in electrical and computer engineering from the University of Illinois at Urbana-Champaign (UIUC), Champaign, IL, USA, in 2010.

He is an Assistant Professor with the School of Electrical and Electronic Engineering, Nanyang Technological University (NTU), Singapore. His research interests include computer vision and machine learning. Particularly, he is focusing on

object recognition, scene analysis, and deep learning. He is an associate editor of *Neurocomputing*.

Prof. Wang was the recipient of the Harriett & Robert Perry Fellowship (2009–2010) and the CS/AI award (2009) while at UIUC.