

Received July 10, 2019, accepted July 26, 2019, date of publication July 30, 2019, date of current version August 16, 2019.

Digital Object Identifier 10.1109/ACCESS.2019.2932080

# RGB-D Scene Recognition via Spatial-Related Multi-Modal Feature Learning

ZHITONG XIONG, YUAN YUAN, (Senior Member, IEEE),  
AND QI WANG<sup>id</sup>, (Senior Member, IEEE)

School of Computer Science and Center for OPTical IMagery Analysis and Learning (OPTIMAL), Northwestern Polytechnical University, Xi'an 710072, China

Corresponding author: Yuan Yuan (y.yuan1.ieee@gmail.com)

This work was supported in part by the National Natural Science Foundation of China under Grant U1864204 and Grant 61773316, in part by the State Key Program of National Natural Science Foundation of China under Grant 61632018, in part by the Natural Science Foundation of Shaanxi Province under Grant 2018KJXX-024, and in part by the Project of Special Zone for National Defense Science and Technology Innovation.

**ABSTRACT** RGB-D image-based scene recognition has achieved significant performance improvement with the development of deep learning methods. While convolutional neural networks can learn high-semantic level features for object recognition, these methods still have limitations for RGB-D scene classification. One limitation is that how to learn better multi-modal features for the RGB-D scene recognition is still an open problem. Another limitation is that the scene images are usually not object-centric and with great spatial variability. Thus, vanilla full-image CNN features maybe not optimal for scene recognition. Considering these problems, in this paper, we propose a compact and effective framework for RGB-D scene recognition. Specifically, we make the following contributions: 1) A novel RGB-D scene recognition framework is proposed to explicitly learn the global modal-specific and local modal-consistent features simultaneously. Different from existing approaches, local CNN features are considered for the learning of modal-consistent representations; 2) key Feature Selection (KFS) module is designed, which can adaptively select important local features from the high-semantic level CNN feature maps. It is more efficient and effective than object detection and dense patch-sampling based methods, and; 3) a triplet correlation loss and a spatial-attention similarity loss are proposed for the training of KFS module. Under the supervision of the proposed loss functions, the network can learn import local features of two modalities with no need for extra annotations. Finally, by concatenating the global and local features together, the proposed framework can achieve new state-of-the-art scene recognition performance on the SUN RGB-D dataset and NYU Depth version 2 (NYUD v2) dataset.

**INDEX TERMS** RGB-D, scene recognition, global and local features, multi-modal feature learning.

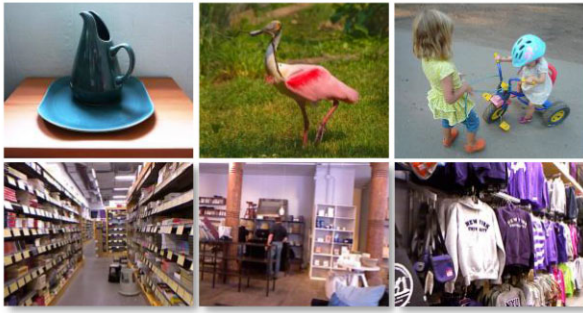
## I. INTRODUCTION

With the advent of deep learning methods especially the convolutional neural networks (CNN), image classification performance has been improved dramatically on the large-scale object-centric image recognition dataset: ImageNet [1]. Although modern CNN architectures such as ResNet [2] can learn more effective representations of image, directly exploiting full-image features is sub-optimal for scene recognition. The reason is that global scene image features cannot capture the great spatial varieties of the scene.

Considering the difference between object recognition and scene classification, varieties of methods [4], [5] have been

proposed for RGB image based scene classification task. Zhou *et al.* [6] released a large scale scene image classification dataset named **Places**, and showed the effectiveness of pre-training CNN parameters on it compared to the ImageNet dataset. To handle the complex geometric variability of scene image, explicitly extracting the object-level or theme-level features has been explored by several methods. The work of [7], [8] and [9] were proposed to leverage the local CNN features for scene classification. These methods firstly extracted features of different scales and locations densely, and then encoded them with the fisher vector (FV) [10]. Although these works can improve the performance with the powerful local features, there exist two obvious disadvantages. One is that merely exploiting the local features neglects the global layout of the scene. Another disadvantage

The associate editor coordinating the review of this manuscript and approving it for publication was Kumaradevan Punithakumar.



**FIGURE 1.** Object-centric images (first row) and scene classification images (second row). Images shown in the first row are selected from ImageNet [1], and the second row images are selected from NYU Depth V2 dataset [3].

is that using densely sampled local features may introduce noise into the final feature encodings, which may further limit the performance.

With the rapid development of depth sensors, RGB-D image based scene classification has attracted increasing research interest. As RGB-D indoor scene images are also not object-centric, several methods [11], [12] were proposed to learn the component-aware semantic features and represent the indoor scene with the combination of object-level features. However, these methods need to accurately detect the objects before scene classification. Thus the performances of these methods deeply rely on the object detection accuracy. Moreover, it is really non-trivial to detect the cluttered objects accurately in the complex indoor scenes.

Although RGB-D scene image can provide extra geometric information compared to the common RGB image, how to learn the multi-modal features effectively is critical for the performance improvement. Many multi-modal representation learning strategies have been proposed to exploit the complementary information of two modalities. The work of Wang *et al.* [13] aimed to minimize the distance of RGB and Depth embeddings. Although enforcing the multi-modal consistency can exclude the noise, it also hinders the modal-complementary feature learning. Li *et al.* [14] proposed a discriminative multi-modal feature learning framework, which learned the distinctive embedding and the correlative embedding simultaneously. However, merely global features are used for multi-modal feature learning and fusion, and local features are neglected in these methods.

Traditional multi-modal feature learning methods usually neglect an important factor: the spatial distribution of features. Based on the fact that depth modality can capture more accurate global scene layout information than RGB modality, we propose to learn modal-specific features from global features and enforce the modality-consistency on selected local features. Since depth modality is also a image, there exists spatial correspondence between the RGB and depth modality. This makes RGB-D image based scene recognition different from visual & text or visual & audio multi-modal feature learning task. However, prior works usually neglect the

spatial distribution of features when extracting modal-distinctive and modal-consistent representations.

To handle the aforementioned issues, in this work, we propose an end-to-end multi-modal feature learning framework, which adaptively selects important local region features and fuses the local and global features together for RGB-D scene recognition. Different from densely patch-sampling based or object detection based approaches, the proposed method selects important local features at different locations on the high-semantic level CNN feature maps. Moreover, we consider the spatial distribution for multi-modal feature learning by encouraging the modality-consistency and modality-correlation on local and global features respectively. Specifically, our contributions can be summarized as follows.

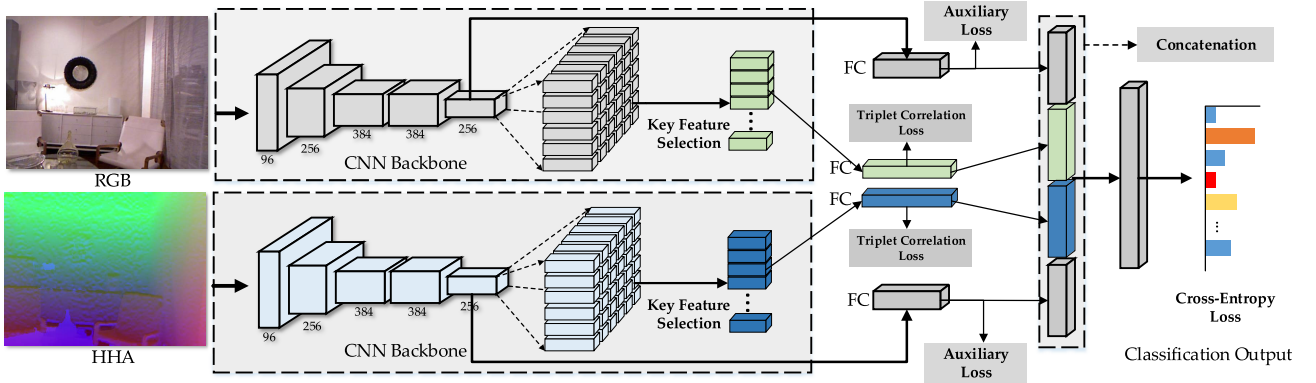
- 1) A novel RGB-D scene recognition framework is proposed to explicitly learn the global modal-specific and local modal-consistent features simultaneously. Different from existing approaches, the spatial distribution of features is considered for multi-modal representation learning.
- 2) Key Feature Selection (KFS) module is designed, which can adaptively select important local features from the high-semantic level CNN feature maps. It is more efficient and effective than object detection and dense patch-sampling based methods.
- 3) A triplet correlation loss and a spatial-attention similarity loss are proposed to learn the local modal-consistent features. With these loss functions, the network can learn the common local patterns between two modalities with no need for extra annotations.

Experiments on two public datasets SUN RGB-D [15] and NYU v2 [3] have shown the effectiveness of the proposed method.

The reminder of this paper is organized as follows. We review the related works in section II. In section III, the details of the proposed method is described. Experimental results and analysis are presented in section IV. Finally, the conclusion is drawn in section V.

## II. RELATED WORK

Many computer vision tasks [16] have achieved great performance improvement with the surge of deep learning methods. However, full-image global CNN features are not flexible enough to represent the complex indoor scene. Thus several local CNN features based methods are proposed for RGB-D scene classification. To learn better local CNN features, Gong *et al.* [7] introduced a multi-scale CNN framework to aggregate densely sampled multi-scale features with the vector of locally aggregated descriptors (VLAD) [17]. The work of [8] and [5] proposed to encode the scene image with multi-scale local activations via the fisher vector (FV) encoding. Song *et al.* [18] firstly trained the model on depth image patches in a weakly-supervised manner, and then fine tuned the model with full image. Nevertheless, the densely



**FIGURE 2.** The whole framework of the proposed method. RGB and depth image are firstly input to two CNN for feature extraction. Then the global modal-specific features are learned by fully connected layers with cross entropy loss. Local modal-consistent features for both RGB and depth modality are learned with the proposed KFS modules. Finally, global and local features are combined together for the final scene recognition.

sampled image patches for feature encoding may contain noise, which decreases the recognition performance.

To handle the aforementioned problems, several methods employ object detection to extract object-level local features. Wang *et al.* [11] attempted to use the CNN region proposals as local features, and combined the local and global features via FV to learn component-aware representations. The work of [12] introduced object detection on RGB-D image to obtain more accurate object-level local features, and they further modeled the object relation among the detected objects. Although improved performance can be achieved, the error accumulation problem of two-stage pipeline methods and higher computational complexity are still limitations.

Multi-modal feature learning strategy is critical for RGB-D scene classification task. To fuse multi-modal features, varieties of strategies have been investigated [19]. Image level multi-modal fusion was proposed in [20] by constructing the RGB-D Laplacian pyramid. Song *et al.* [15] fused the two modal features by concatenating two-stream CNN features to one fully connected layer. The work of [21] employed a three-stream CNN to combine RGB branch and two depth modal features by using element-wise summation. To learn modal-consistent features, Wang *et al.* [22] enforced the network to learn common features between RGB and depth images. Li *et al.* [14] aimed to learn the correlative and distinctive embeddings between the two modalities simultaneously. However, enforcing modal consistency hinders the complementary feature learning, which may decrease the performance. Moreover, these multi-modal learning methods do not take local features into consideration.

### III. OUR METHOD

The whole proposed framework is shown in Fig. 2. Firstly, three RGB and depth (HHA encoded [23]) image pairs are sampled to input the network as a training triplet. After the feature extraction through a two-branch CNN, global modal-specific features are learned with the FC (Fully Connected) layers and auxiliary loss. Meanwhile, local important

region features are selected with the KFS module. Then the global and local features of two modalities are concatenated together for the final scene classification. The details of the proposed framework will be described in the following sections.

#### A. KEY LOCAL FEATURE SELECTION

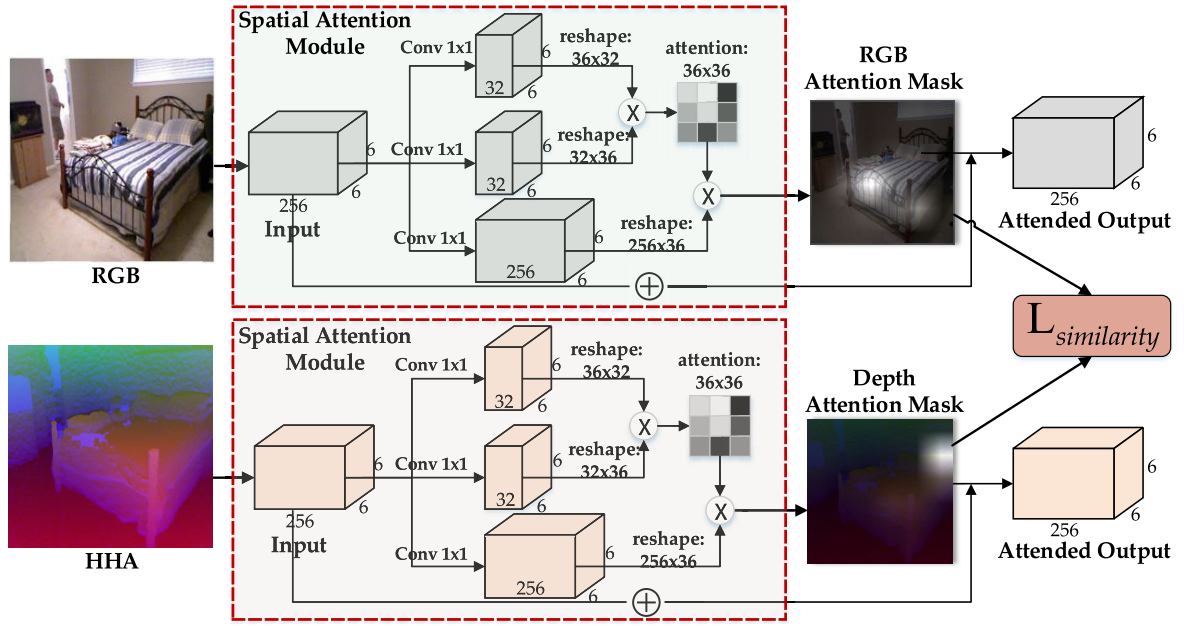
The great intra-class variation RGB-D indoor scene makes the classification challenging. From the second row of Fig. 1 we can see that, there are three dramatically different images with the same “book store” class label. Considering this, we employ local object-level features to reduce the intra-class variation.

Different from patch-based and object detection based methods, in this work, we aim to select key region features from the high semantic level CNN feature maps. As CNN learns local features by convolution operation, the local spatial context information is embedded into the feature vectors of the final feature maps. Suppose that  $F_{rgb}$  is the final feature of RGB branch for one training sample in the input triplet.

As the scene image can usually be represented by several typical objects or themes, we opt to select  $K$  local object or theme-level features from  $F_{rgb} \in \mathbb{R}^{(N,C,H,W)}$  for classification.  $N$  is the batch size. For feature selection, it is critical to define the criteria to measure the importance of features. To learn which features are more important, we employ the spatial attention based models to enhance the local feature selection module in this work. Specifically, the non-local networks are employed as the spatial attention module, which can be formulated as follows.

$$\begin{aligned} A_{rgb} &= \text{softmax}(\theta(F_{rgb})^T \phi(F_{rgb})), \\ F'_{rgb} &= A_{rgb} g(F_{rgb}), \end{aligned} \quad (1)$$

where  $g$  is a  $1 \times 1$  convolutional layer with the same output channel number as the input features.  $\theta$  and  $\phi$  are  $1 \times 1$  convolutional layers for transforming the input feature  $F_{rgb}$  in non-local networks.  $\theta$ ,  $\phi$  and  $g$  are convolutional layers for learning the attention mask, which are different for RGB



**FIGURE 3.** The illustration for the attention mask similarity loss of the proposed KFS module. The attention masks of RGB and depth image are encouraged to be similar to learn the local modal-consistent features.

and depth modality. In this work, the dot-product similarity is used to measure the similarity between features at different spatial positions. Then the final spatial attention results are obtained by

$$FA_{rgb} = F_{rgb} + \gamma F'_{rgb}, \quad (2)$$

where  $\gamma$  is a learnable parameter, and its initial value is set to 0. Intuitively, features with higher response should be more important than the lower ones. Thus we sum over all the channels of  $FA_{rgb}$  to get a response map  $F_{resp} \in \mathbb{R}^{(N,H,W)}$ , and  $H, W$  are height and width of the response map.

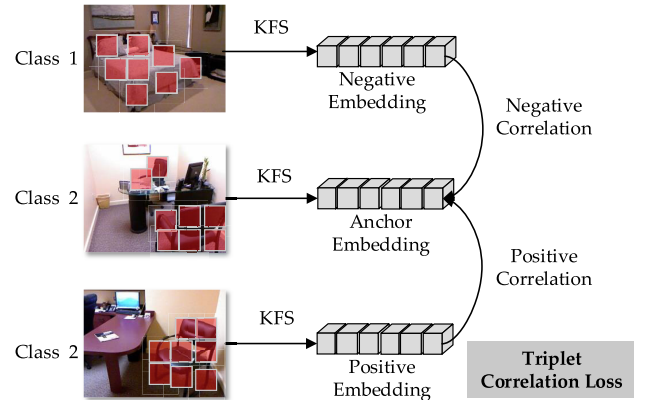
Then we reshape the response map to  $\{N, H \times W\}$ , and sort them to find the  $K$  highest response indexes. With the ‘Sort and Select’ strategy,  $K$  local feature vectors are selected to represent the scene.

However, merely considering the filter response is insufficient to select discriminative features for scene classification. Since we aim to select local features which are critical for classifying the scene of different classes, the triplet correlation loss is proposed to regularize the local feature selection process.

The triplet loss module is shown in Fig.4. The selected local features of each sample in the input triplet can be denoted as  $E_p, E_a, E_n \in \mathbb{R}^{N,C \times K}$ , where  $E_p$  and  $E_a$  are positive and anchor features with the same class label, and  $E_n$  is the negative one with different label. The triplet correlation loss can be formulated as

$$L_{rgb\_trip\_corr} = \max\{\rho(E_a, E_p) - \rho(E_a, E_n) + \alpha, 0\},$$

$$\rho(x, y) = \frac{\langle x, y \rangle}{\|x\| \|y\|}, \quad (3)$$



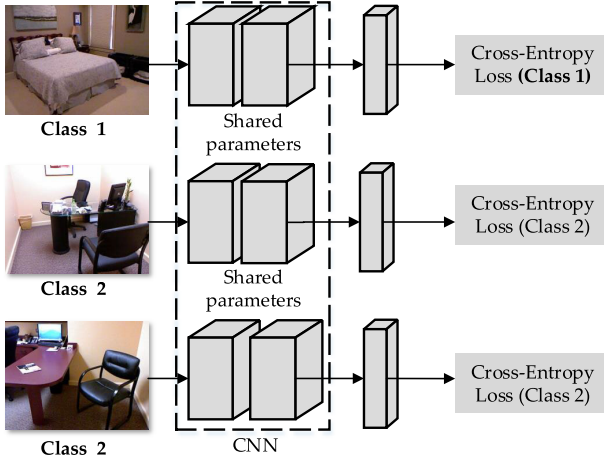
**FIGURE 4.** The illustration of local feature selection and learning module. We aim to select local features which are correlated between different images from the same scene class.

where  $L_{rgb\_trip\_corr}$  is the triplet correlation loss for local feature selection of the RGB modality, and the loss computation for depth modality is similar to the RGB modality.

## B. ENFORCING THE MULTI-MODAL FEATURE CONSISTENCY

RGB-D image based feature learning is quite different from visual/audio or visual/text multi-modal feature learning task. Since RGB and depth modality are both images and they are spatially aligned, we can further enhance the modality consistency using the attention mask depicted in section III-A. As aforementioned, the spatial attention module is designed to capture the local part features without the need for extra annotations.





**FIGURE 5.** The illustration of global modal-specific feature learning module. For each modality input (a triplet in this work), the modal-specific features are learned by separately training with the cross-entropy loss.

The experimental results of the Key Local Feature selection module reveal that the spatial attention map for RGB and depth modality are surprisingly similar, as shown in Fig. 6. Inspired by this observation, we further design a loss term to enforce the consistency of multiple modality attention maps. The detailed illustration of this loss is shown in Fig. 3. The spatial attention module are used to allow the model to focus on key local features for both the RGB and depth input images. Thus two attention masks are obtained for the two modalities. Based on the observation that attention masks for RGB and depth modality are similar in spatial distribution. We further propose a loss term to maximize the similarity between two attention masks of different modalities.

Specifically, suppose the attention maps of RGB and depth modality are  $A_{rgb}$  and  $A_d$  respectively, the similarity loss  $L_{sim}$  can be computed as

$$L_{sim} = \frac{1}{2} \|A_{rgb} - A_d\|_2^2. \quad (4)$$

By encouraging the network to focus on features at similar spatial positions, the proposed framework can learn more representative modal-consistent features.

### C. DISCRIMINATIVE MULTI-MODAL GLOBAL FEATURE LEARNING

As global features are useful for describing the scene layouts, they are important for scene classification. To make full use of each modality, the modal-specific global features are extracted for RGB and depth modality respectively.

Specifically, we first sample three images as a triplet input, which consists of two images with the same class label and one image with different class labels. For simplicity, the triplet samples are denoted as  $\{x_1, x_2, x_3\}$  and their labels as  $\{y_1, y_2, y_3\}$ . In this triplet, we set  $y_1 = y_2$ . For these three samples, cross-entropy loss is used for image classification.

As the global feature learning process for RGB and depth modality are similar, we take the RGB branch for example. For simplicity, we represent the CNN feature learning as

$$F_{rgb} = f_{rgb}(x), \quad (5)$$

and the three learned global embeddings  $G_p, G_a, G_n$  are obtained by a fully connected layer.  $G_p$  is the feature of positive sample  $y_1$ , which has the same class label with the anchor sample  $y_2$ .  $G_a$  is the feature of  $y_2$ , and it has different class label with negative sample  $y_3$ , whose feature embedding is  $G_n$ .

To learn the modal-specific features, we propose to train the two branches of CNN separately. Since RGB and depth modality data contains different information, training two branches of CNN separately can force the CNN to learn specific features for each modality. In this work, cross-entropy loss is applied for training the two branches of CNN separately. This can be represented as

$$L_{aux}^{rgb} = \sum_{i=1}^3 L_{CE}(\hat{y}_i, y_i), \quad (6)$$

where  $L_{aux}^{rgb}$  is the auxiliary loss function for discriminative feature learning.  $\hat{y}_i$  and  $y_i$  are the class predictions of the global features and ground truth respectively. The global feature learning for RGB modality is illustrated in Fig. 5.

For the depth modality, the loss computation is similar to the RGB modality depicted above.

### D. MULTI-MODAL GLOBAL AND LOCAL FEATURE FUSION

To fuse the learned global and local features together for RGB and depth modality, we concatenate them into a multi-modal global and local feature vector for the final scene classification. This can be denoted as

$$F_{mmgl} = \text{concat}(E_{rgb}, E_d, G_{rgb}, G_d), \quad (7)$$

where  $F_{mmgl}$  is the multi-modal global and local feature vector.  $E_{rgb}$  and  $E_d$  are the selected local features for RGB and depth modality.  $G_{rgb}$  and  $G_d$  are the global features of the RGB and depth modality.

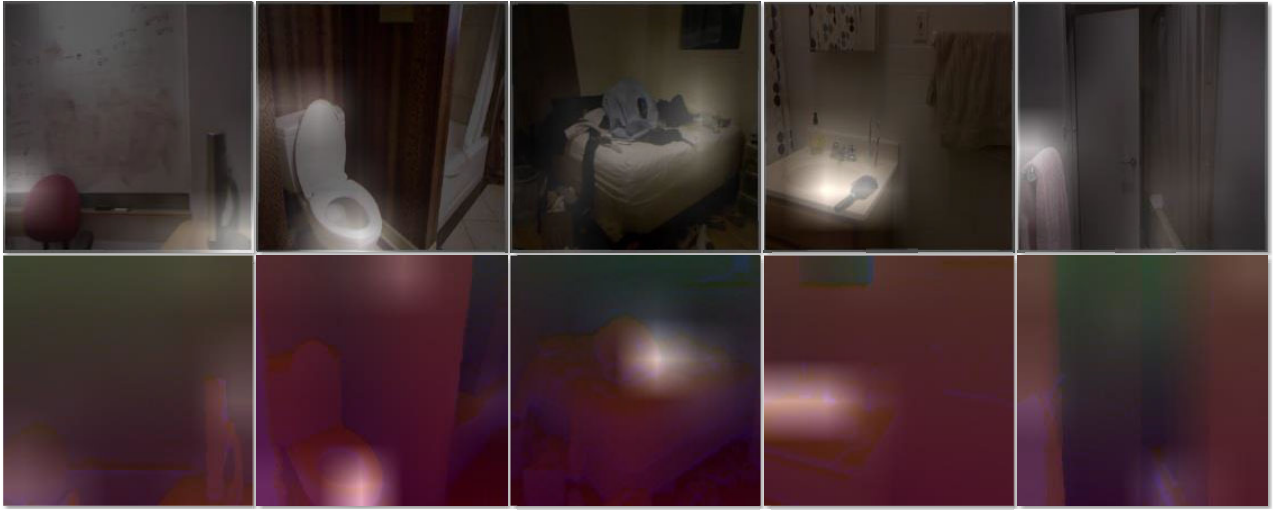
After passing  $F_{mmgl}$  through a fully connected layer, the final classification result can be predicted with an extra softmax layer. For all the three samples in the input triplet, cross entropy is used as the final classification loss, which can be represent as

$$L_{cls} = \sum_{i=1}^3 L_{CE}(\hat{y}_i, y_i), \quad (8)$$

where  $\hat{y}_i$  and  $y_i$  are the class predictions of the final multi-modal features and ground truth respectively.

Finally, the overall loss of the proposed framework consists of two global modal-specific auxiliary loss functions, two triplet correlation loss functions and one final classification loss function. Thus the total loss function can be formulated as

$$L = L_{cls} + \lambda_1 L_{aux} + \lambda_2 L_{Trip\_corr} + \lambda_3 L_{sim}, \quad (9)$$



**FIGURE 6.** The illustration of the attention map for selecting key local features. As shown in this figure, the attention masks for RGB and depth modality are similar in spatial distribution, which indicates that the local feature learning processes for different modality share the same pattern. Based on this observation, we further proposed an attention mask similarity loss to enhance the local modal-consistent feature learning.

where the term  $L_{aux}$  consists of RGB and depth auxiliary loss.  $L_{aux}$  is computed by

$$L_{aug} = L_{aux}^{rgb} + L_{aux}^d. \quad (10)$$

The triplet correlation loss also includes two terms for both RGB and depth modality, which is defined as

$$L_{Trip\_corr} = L_{rgb\_trip\_corr} + L_{d\_trip\_corr}. \quad (11)$$

The proposed framework takes triplet as input and learns the global and local multi-modal features with the auxiliary loss and triplet correlation loss described above. The whole framework can be trained in an end-to-end manner and can be easily implemented using modern deep learning frameworks.

#### IV. EXPERIMENTS

We evaluate the proposed method on two public datasets: SUN RGB-D [15] and NYU Depth Dataset version 2 [3]. There are 10,355 RGB images with corresponding depth images in SUN RGB-D dataset, and they are divided into 19 categories. Following the previous experimental settings [15], we use 4,845 images for training, 4,659 images for testing. NYUD v2 contains 1449 images, and they are divided into 10 categories including 9 common indoor scene types and one ‘others’ category. In this dataset, 795 images are used for training and 654 for testing following the setting in [24].

Mean-class accuracy is used as the evaluation measurement in this work to compare with previous methods. It is computed by averaging the precisions for all the categories, i.e., the diagonal elements of the confusion matrix. The mean-class accuracy can be defined as follows.

$$MeanAcc = \frac{1}{C} \sum_{c=1}^C \frac{correct_c}{Num_c}, \quad (12)$$

where  $correct_c$  is the number of correctly predicted samples of class  $c$ , and  $Num_c$  is the total number of samples of class  $c$ .

#### A. PARAMETERS SETUP

We compute the HHA encodings with the released code from [23]. As the training samples are scarce for deep CNN, data augmentation is used in our work. The input images in triplet are firstly resized to  $224 \times 224$ , and then random horizontal flip and random erasing [25] are used for each modality at a probability of 0.5. To compare it with existing work, AlexNet [26] with pre-trained parameters on Places dataset is used as the back-bone network. Adam [27] optimizer is employed with a initial learning rate of  $1e-4$ . The learning rate is reduced by a fraction of 0.9 every 80 epochs during training. The batch size is set to 64 with shuffle and 300 epochs are used to train the proposed framework. For the multi-task training, we set the parameters  $\lambda_1$  and  $\lambda_2$  to 1 in all of our experiments. The parameter  $\lambda_3$  is set to 0.001 in all the experiments.

#### B. SUN RGB-D DATASET

We compare six state-of-the-art methods on SUN RGB-D dataset. Among them, Song et al. [15] took RGB and HHA encoding as input for scene classification. [28] combined the scene recognition and semantic segmentation tasks into one multi-task framework. Zhu et al. [29] considered the intra-class and inter-class correlations for scene classification. Wang et al. [11] and Song et al. [12] introduced object detection based local feature learning methods. [14] proposed a framework to learn distinctive and correlative features simultaneously. From the results in Table 1, our method achieves state-of-the-art performance 55.9%, which is better than object detection based methods [12], [11].

Among these compared methods, the work of Li et al. [14] is the most related work with ours. They achieved state-of-the-art performance by learning the modal-distinctive and modal-correlative features simultaneously. However, local features are not considered in their work, which is important

**TABLE 1.** Experimental results on SUN RGB-D dataset.

	Methods	Local Features	Accuracy(%)
State-of-the-art	Song et al. [15]	No	39.0 %
	Liao et al. [28]	No	41.3%
	Zhu et al. [29]	No	41.5%
	Wang et al. [11]	CNN Proposals	48.1%
	Song et al. [18]	Local Patches	52.4%
	Song et al. [12]	Object Detection	54.0%
	Li et al. [14]	No	54.6%
Proposed	Our method	Key Feature Selection	<b>55.9%</b>

**TABLE 2.** Experimental results on NYUD v2 dataset.

	Methods	Local Features	Accuracy(%)
State-of-the-art	Gupta et al. [30]	No	45.4 %
	Wang et al. [11]	CNN Proposals	63.9%
	Li et al. [14]	No	65.4%
	Song et al. [18]	Local Patches	65.8%
	Song et al. [12]	Object Detection	66.9%
Proposed	Our method	Key Feature Selection	<b>67.8%</b>

for scene recognition tasks. By extracting and combining the global and local features, the proposed method in this work can obtain better performance, which indicates the effectiveness of the designed KFS module.

For local feature extraction, Song *et al.* [12] proposed to learn local object-level features by performing the object detection task in advance of scene recognition. They further took the relationships between objects into consideration and achieved quite good performance. However, although global and local features are considered in their work, they neglect the modal-correlation and modal-distinction for multi-modal feature learning. By integrating the global/local and modal-specific/consistent feature learning processes, the proposed framework can achieve better scene recognition performance.

### C. NYUD V2 DATASET

On NYU v2 dataset, five state-of-the-art methods are compared. Among them, Song *et al.* [18] tried to learn depth features by firstly training the network on local depth patches. As presented in Table 2, our approach achieves better performance (accuracy 67.8%) than state-of-the-art method (66.9%) on NYUD v2 dataset. Moreover, it is worth mentioning that no object detection is needed in our approach. In general, the experimental results on NYUD v2 dataset is similar to SUN RGB-D dataset. The proposed method can obtain better recognition performance than existing state-of-the-art methods, which indicates the effectiveness of our framework.

Additionally, ablation study on NYU v2 dataset is conducted for more comprehensive evaluations of the proposed method. As shown in Table 3, merely employing single modality information can only achieve limited performance. With the discriminative global feature learning module, ‘RGB-D Global(Discriminative Learning)’ can obtain 64.1% accuracy, which improves 2.6% than baseline 61.5%. We also evaluate the effect of the triplet loss of KFS module. The results show that 65.3% accuracy can be obtained without

**TABLE 3.** Ablation study on NYUD v2 dataset.

Methods	Accuracy(%)
RGB (Single Modality)	53.5%
Depth (Single Modality)	51.1%
RGB-D (Multi-modal Baseline)	61.5%
RGB-D Global (Discriminative Learning)	64.1%
RGB-D Global & Local (KFS, w/o Triplet Loss)	65.3%
RGB-D Global & Local (KFS, w/o $L_{sim}$ )	66.5%
RGB-D Global & Local (KFS, Full Model)	<b>67.8%</b>

the triplet loss term, which is worse than the performance of KFS module with triplet loss (66.5%). Additionally, to show the effect of the attention mask similarity loss  $L_{sim}$ , we do experiment without the loss term  $L_{sim}$ , and 66.5% mean-class accuracy is obtained. The comparison results indicate the effectiveness of the proposed sub-modules.

To validate the effect of the proposed KFS module, we compared the ‘‘RGB-D (Multi-modal Baseline)’’ method with ‘‘RGB-D Global & Local (KFS, w/o Triplet Loss)’’. As the RGB-D baseline merely employs the global features, it achieves a lower performance of 61.5%. By exploiting the local features with KFS module, the ‘‘RGB-D Global & Local (KFS, w/o Triplet Loss)’’ achieves a higher performance of 65.3%. This comparison indicates that the proposed KFS module can learn effective local representations which are complementary to the global features. For the validation of triplet correlation loss, we compare the ‘‘RGB-D Global & Local (KFS, w/o Triplet Loss)’’ with ‘‘RGB-D Global & Local (KFS, w/o  $L_{sim}$ )’’. Since ‘‘RGB-D Global & Local (KFS, w/o  $L_{sim}$ )’’ method uses the triplet correlation loss, it can obtain a 66.5% mean-class accuracy, which is better than ‘‘RGB-D Global & Local (w/o Triplet Loss)’’ method. This comparison verifies the effectiveness of the proposed triplet correlation loss.

Some examples of the selected key regions of the proposed method is presented in Fig. 7. From the figure we can see that some common object-level features are selected for the same



**FIGURE 7.** The illustration of the attention map for selecting key local features. The upper row shows the attention maps for RGB modality, and the lower row shows the attention map for depth modality. It is worth mentioning that we still use RGB image instead of its HHA encoding for more clear presentation. We can see that the attention masks for RGB and depth modality have similar spatial distributions for learning modal-consistent local features.

scene class. In Fig. 7, the upper row shows the attention maps for RGB modality, and the lower row shows the attention map for depth modality. It is worth mentioning that we still use RGB image instead of its HHA encoding for more clear presentation.

#### D. DISCUSSION

To sum up, experiments on public datasets have indicated the effectiveness of the proposed method. From the experimental results, we find that local features, i.e., CNN intermediate features are critical to the performance improvement of scene classification. By combining the local CNN features with the global features, the accuracy can be boosted. This reveals that the selected local features are complementary to global features for scene classification task. Moreover, the proposed local feature selection module can be trained jointly with scene recognition task in an end-to-end manner, which is more efficient than object-detection based method. Additionally, traditional modal-consistent representations are extracted with global features, while the experiments of this work indicate that local modal-consistent features can be useful for RGB-D scene recognition.

From the experiments we observe that foreground objects share similar pattern between two modalities, and the performance can be boosted by further enhancing the pattern

similarity (similarity between attention masks). However, depth modality contains more information about the global scene layout, the global background features may be more suitable for learning the modal-distinctive features. Enforcing the dissimilarity on selected local features of two modalities cannot improve the performance. Considering that the KFS module is proposed to select foreground object-level features, the similarity loss is more useful for KFS module.

The noise in depth images can do harm to the learned features in RGB-D image based tasks. For depth images in the RGB-D scene recognition datasets, the noise usually appears on the background areas (away from image center). However, the proposed method mainly focuses on learning the features of foreground objects by KFS module. Since the foreground objects usually contain no noise, learning local foreground features with the KFS module can help to alleviate the effect of the depth noise.

#### V. CONCLUSION

In this paper, we propose a compact and effective framework for RGB-D scene recognition, which fuses local and global features together to improve the recognition performance. To extract effective local features of key objects or themes, we propose a key feature selection (KFS) module, which adaptively selects key local features under the supervision



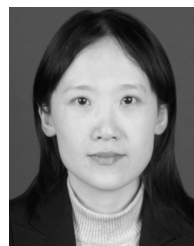
of a triplet correlation loss and a multi-modal consistency loss. With this module, the proposed method can learn more discriminative local representations. Besides, global modal-specific features are extracted for the two modalities respectively under the supervision of the proposed auxiliary loss. By concatenating the global and local features, the proposed framework can achieve new state-of-the-art scene recognition performance on the SUN RGB-D dataset and NYU Depth version 2 (NYUD v2) dataset.

## REFERENCES

- [1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. CVPR*, Jun. 2009, pp. 248–255.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2016, pp. 770–778. doi: 10.1109/CVPR.2016.90.
- [3] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from RGBD images," in *Proc. ECCV*, 2012, pp. 746–760.
- [4] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 2, 2006, pp. 2169–2178.
- [5] M. Dixit, S. Chen, D. Gao, N. Rasiwasia, and N. Vasconcelos, "Scene classification with semantic Fisher vectors," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2015, pp. 2974–2983.
- [6] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, "Learning deep features for scene recognition using places database," in *Proc. Adv. Neural Inf. Process. Syst.*, 2014, pp. 487–495.
- [7] Y. Gong, L. Wang, R. Guo, and S. Lazebnik, "Multi-scale orderless pooling of deep convolutional activation features," in *Proc. ECCV*, 2014, pp. 392–407.
- [8] D. Yoo, S. Park, J.-Y. Lee, and I.-S. Kweon, "Fisher kernel for deep neural activations," *CoRR*, vol. abs/1412.1628, 2014.
- [9] Z. Zuo, G. Wang, B. Shuai, L. Zhao, Q. Yang, and X. Jiang, "Learning discriminative and shareable features for scene classification," in *Proc. ECCV*, 2014, pp. 552–568.
- [10] J. Sánchez, F. Perronnin, T. Mensink, and J. J. Verbeek, "Image classification with the Fisher vector: Theory and practice," *Int. J. Comput. Vis.*, vol. 105, no. 3, pp. 222–245.
- [11] A. Wang, J. Cai, J. Lu, and T.-J. Cham, "Modality and component aware feature fusion for RGB-D scene classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2016, pp. 5995–6004.
- [12] X. Song, C. Chen, and S. Jiang, "RGB-D scene recognition with object-to-object relation," in *Proc. ACM Multimedia Conf.*, 2017, pp. 600–608.
- [13] A. Wang, J. Lu, J. Cai, T.-J. Cham, and G. Wang, "Large-margin multi-modal deep learning for RGB-D object recognition," *IEEE Trans. Multimedia*, vol. 17, no. 11, pp. 1887–1898, Nov. 2015.
- [14] Y. Li, J. Zhang, Y. Cheng, K. Huang, and T. Tan, "DF<sup>2</sup>Net: Discriminative feature learning and fusion network for RGB-D indoor scene classification," in *Proc. 32nd AAAI Conf. Artif. Intell.*, 2018, pp. 7041–7048.
- [15] S. Song, S. P. Lichtenberg, and J. Xiao, "Sun RGB-D: A RGB-D scene understanding benchmark suite," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2015, pp. 567–576.
- [16] X. Li, Z. Yuan, and Q. Wang, "Unsupervised deep noise modeling for hyperspectral image change detection," *Remote Sens.*, vol. 11, no. 3, p. 258, 2019.
- [17] H. Jegou, M. Douze, C. Schmid, and P. Pérez, "Aggregating local descriptors into a compact image representation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 3304–3311.
- [18] X. Song, L. Herranz, and S. Jiang, "Depth CNNs for RGB-D scene recognition: Learning from scratch better than transferring from RGB-CNNs," in *Proc. AAAI*, 2017, pp. 4271–4277.
- [19] Q. Wang, M. Chen, F. Nie, and X. Li, "Detecting coherent groups in crowd scenes by multiview clustering," *IEEE Trans. Pattern Anal. Mach. Intell.*, to be published.
- [20] C. Couprie, C. Farabet, L. Najman, and Y. LeCun, "Indoor semantic segmentation using depth information," 2013, *arXiv:1301.3572*. [Online]. Available: <https://arxiv.org/abs/1301.3572>
- [21] X. Song, S. Jiang, and L. Herranz, "Combining models from multiple sources for RGB-D scene recognition," in *Proc. IJCAI*, 2017, pp. 4523–4529.
- [22] A. Wang, J. Cai, J. Lu, and T.-J. Cham, "MMSS: Multi-modal sharable and specific feature learning for RGB-D object recognition," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2015, pp. 1125–1133.
- [23] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, "Learning rich features from RGB-D images for object detection and segmentation," in *Proc. ECCV*. Springer, 2014, pp. 345–360.
- [24] S. Gupta, P. Arbeláez, and J. Malik, "Perceptual organization and recognition of indoor scenes from RGB-D images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2013, pp. 564–571.
- [25] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," 2017, *arXiv:1708.04896*. [Online]. Available: <https://arxiv.org/abs/1708.04896>
- [26] 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.
- [27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, *arXiv:1412.6980*. [Online]. Available: <https://arxiv.org/abs/1412.6980>
- [28] Y. Liao, S. Kodagoda, Y. Wang, L. Shi, and Y. Liu, "Understand scene categories by objects: A semantic regularized scene classifier using convolutional neural networks," in *Proc. ICRA*, May 2016, pp. 2318–2325.
- [29] H. Zhu, J.-B. Weibel, and S. Lu, "Discriminative multi-modal feature fusion for RGBD indoor scene recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2016, pp. 2969–2976.
- [30] S. Gupta, P. Arbeláez, R. Girshick, and J. Malik, "Indoor scene understanding with RGB-D images: Bottom-up segmentation, object detection and semantic segmentation," *Int. J. Comput. Vis.*, vol. 112, no. 2, pp. 133–149, 2015.



**ZHITONG XIONG** received the M.E. degree from the Northwestern Polytechnical University, Xi'an, China, where he is currently pursuing the Ph.D. degree with the School of Computer Science and the Center for OPTical Imagery Analysis and Learning (OPTIMAL). His research interests include computer vision and machine learning.



**YUAN YUAN** (M'05–SM'09) is currently a Full Professor with the School of Computer Science and the Center for Optical Imagery Analysis and Learning, Northwestern Polytechnical University, Xi'an, China. She has authored or coauthored over 150 papers, including about 100 in reputable journals, such as the IEEE TRANSACTIONS AND PATTERN RECOGNITION, as well as the conference papers in CVPR, BMVC, ICIP, and ICASSP. Her current research interests include visual information processing and image/video content analysis.



**QI WANG** (M'15–SM'15) received the B.E. degree in automation and the Ph.D. degree in pattern recognition and intelligent systems from the University of Science and Technology of China, Hefei, China, in 2005 and 2010, respectively. He is currently a Professor with the School of Computer Science and the Center for Optical Imagery Analysis and Learning, Northwestern Polytechnical University, Xi'an, China. His research interests include computer vision and pattern recognition.

• • •