# Scene Image Representation by Foreground, Background and Hybrid Features

Chiranjibi Sitaula*, Yong Xiang, Sunil Aryal, Xuequan Lu

*School of Information Technology, Deakin University, Geelong, Victoria 3216, Australia*

**Abstract**

Previous methods for representing scene images based on deep learning primarily consider either the foreground or background information as the discriminating clues for the classification task. However, scene images also require additional information (hybrid) to cope with the inter-class similarity and intra-class variation problems. In this paper, we propose to use hybrid features in addition to foreground and background features to represent scene images. We suppose that these three types of information could jointly help to represent scene image more accurately. To this end, we adopt three VGG-16 architectures pre-trained on ImageNet, Places, and Hybrid (both ImageNet and Places) datasets for the corresponding extraction of foreground, background and hybrid information. All these three types of deep features are further aggregated to achieve our final features for the representation of scene images. Extensive experiments on two large benchmark scene datasets (MIT-67 and SUN-397) show that our method produces the state-of-the-art classification performance.

*Keywords:* Image processing, Machine learning, Classification, Deep learning, Image representation, Hybrid features.

---

*Corresponding Author

 *Email address:* csitaul@deakin.edu.au (Chiranjibi Sitaula )

## 1. Introduction

With the prevalent and rising demand of robotics and video surveillance, image representation has been a very important (Sitaula et al., 2019c) field to improve classification and recognition accuracies . However, the image representation depends on the problem domain, because we need to represent the images according to contents present in the images and all images can hardly be represented by a single features extraction method. Broadly, image representation methods are categorized into two categories: content-based image representation methods and context-based image representation methods.

Content-based image representation methods (Zeglazi et al., 2016; Oliva, 2005; Oliva & Torralba, 2001; Dalal & Triggs, 2005; Lazebnik et al., 2006; Wu & Rehg, 2011; Xiao et al., 2014; Margolin et al., 2014; Quattoni & Torralba, 2009; Zhu et al., 2010; Li et al., 2010; Parizi et al., 2012; Juneja et al., 2013; Lin et al., 2014; ShenghuaGao & Liang-TienChia, 2010; Perronnin et al., 2010; Xiao et al., 2010; Sánchez et al., 2013; Gong et al., 2014; Kuzborskij et al., 2016; He et al., 2016; Zhou et al., 2016; Zhang et al., 2017; Tang et al., 2017; Guo et al., 2018; Guo & Lew, 2016; Bai et al., 2019; Zhou et al., 2014; Simonyan & Zisserman, 2014; Yang & Ramanan, 2015; Wu et al., 2015; Dixit et al., 2015) rely on the visual content information of the scene images. These features are either based on traditional computer vision based algorithms (Zeglazi et al., 2016; Oliva, 2005; Oliva & Torralba, 2001; Dalal & Triggs, 2005; Lazebnik et al., 2006; Wu & Rehg, 2011; Xiao et al., 2014; Margolin et al., 2014; Quattoni & Torralba, 2009; Zhu et al., 2010; Li et al., 2010; Parizi et al., 2012; Juneja et al., 2013; Lin et al., 2014; ShenghuaGao & Liang-TienChia, 2010; Perronnin et al., 2010; Xiao et al., 2010; Sánchez et al., 2013) or deep learning based algorithms (Gong et al., 2014; Kuzborskij et al., 2016; He et al., 2016; Zhou et al., 2016; Zhang et al., 2017; Tang et al., 2017; Guo et al., 2018; Guo & Lew, 2016; Bai et al., 2019; Zhou et al., 2014; Simonyan & Zisserman, 2014; Yang & Ramanan, 2015; Wu et al., 2015; Dixit et al., 2015). Traditional computer vision based algorithms are more suitable for specific types of images such as texture images. However,
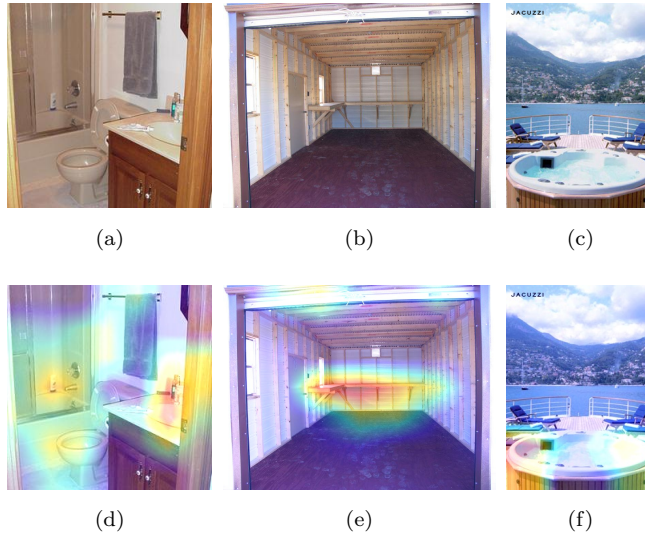
Figure 1: First row shows the scene images belonging to three different categories ((a) toilet, (b) garage, and (c) bathtub outdoor) and the second row shows their corresponding discriminating clues ((d) foreground information for toilet, (e) background information for garage, and (c) hybrid information for bathtub outdoor)). Note that all the feature maps are extracted from the $5^{th}$ pooling layer of VGG-16 models, which are pre-trained on ImageNet, Places, and hybrid datasets (ImageNet+Places).

the recent studies have shown that deep learning based algorithms have higher classification performance than the traditional computer vision based algorithms for the complex scene images involving objects and their associations.

Context-based image representation (Zhang et al., 2017; Wang & Mao, 2019; Sitaula et al., 2019b) addresses the difficult problem of representing the ambiguous images including between-class similarity and within-class dissimilarity. These works are mostly performed based on the exploitation of human annotations/descriptions, with regard to the similar scene images of the input query image on the web. Nevertheless, web crawling and features extraction based on such approaches could be sometimes impractical due to the labor-intensive computational requirements and multiple levels of pre-processing such as tokenization of raw crawled texts, language translation, stemming and lematization,

(a) Foreground          (b) Background
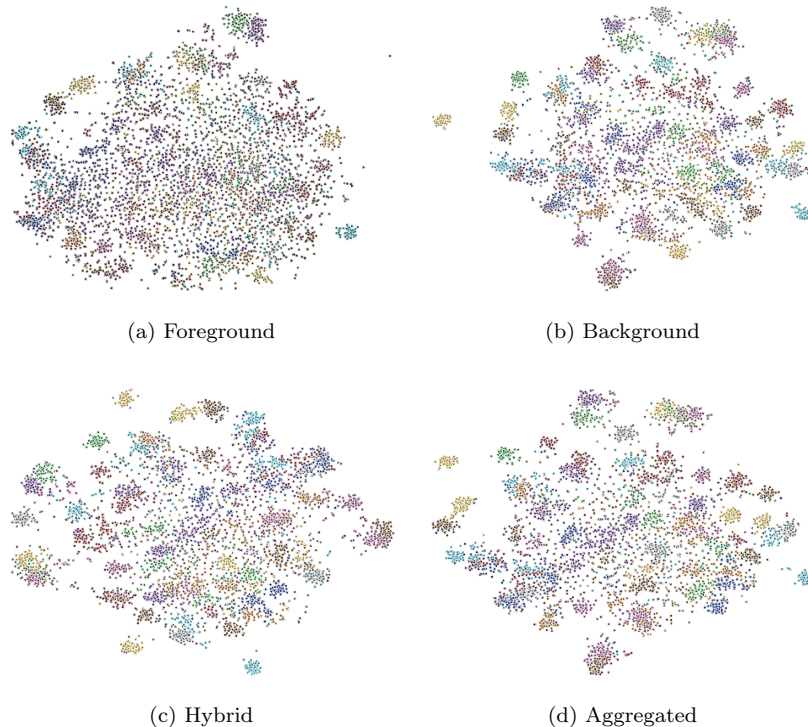
(c) Hybrid          (d) Aggregated

Figure 2: The t-SNE visualization of scatter plots showing the discriminability of features based on the (a) foreground, (b) background, (c) hybrid, and (d) aggregated information. Note that multiple colors represent categories (67 categories) in the MIT-67 dataset.

etc.

The existing methods in the literature primarily focus on either foreground or background information , which may not be sufficient for accurate representation of varying types of scene images. First, different types of scene images may require different types of information to distinguish them accurately. Fig. 1 shows an example. In the figure, three types of scene images requires three different types of information for their better separability. Second, scene images usually involve inter-class similarity and intra-class dissimilarity issues. It may require additional information (hybrid) to foreground and background information to improve the separability.

To bridge the aforementioned gaps above, we perform the fusion of three different types of information including foreground, background, and hybrid for each image. For this, we extract the foreground, background, and hybrid information of each image with the help of VGG-16 models (Simonyan & Zisserman, 2014) pre-trained on ImageNet (Deng et al., 2009), Places (Zhou et al., 2017), and both (ImageNet + Places), respectively. We choose the VGG-16 model due to its simple architecture yet prominent features extraction capability (Guo & Lew, 2016; Guo et al., 2018; Bai et al., 2019). To achieve the corresponding features from the VGG-16 models, we utilize a higher level pooling layer ($p\_5$) (Guo et al., 2018; Guo & Lew, 2016) as the features extraction layer, because we found that the $p\_5$ layer yields highly separable features than other layers (see detail in Section 4.4) . Finally, we aggregate these three types of features to achieve our final features for the classification. The separability of our aggregated features (combined) and individual features are shown in Fig. 2, using t-SNE (t-Distributed Stochastic Neighbor Embedding) scatter plot visualization tool.

The main **contributions** of this paper are summarized as follows.

(a) We propose a novel method for image representation by identifying three different types of information (foreground, background, and hybrid) and fusing them.

(b) We design an effective scheme to aggregate three important types of information using three different pre-trained deep learning models (VGG-16). We analyze five pooling layers of VGG-16 and choose the best features extractor in this work.

(c) We evaluate our method on two common scene benchmark datasets: MIT-67 (Quattoni & Torralba, 2009) and SUN-397 (Xiao et al., 2010). Evaluation results show that our method produces features with better separability and results in the state-of-the-art classification performance.

The rest of the paper is organized as follows. Section 2 reviews the previous image representation methods for the scene images. Section 3 elaborates our

proposed method in a step-wise manner. Section 4 explains the implementation details, the comparisons with the previous methods and the ablative studies. Section 5 concludes this work.

## 2. Related work

In general, image representations can be divided into two types: content-based and context-based.

### 2.1. Content-based image representation

Content-based image representation methods are further categorized into two subgroups: traditional computer-vision based algorithms and deep-learning based algorithms. Traditional computer vision-based algorithms primarily depend on the traditional feature extraction methods such as GIST-color (Oliva & Torralba, 2001), Generalized Search Trees (GIST) (Oliva, 2005), Histogram of Gradient (HOG) (Dalal & Triggs, 2005), Spatial Pyramid Matching (SPM) (Lazebnik et al., 2006), RoI (regions of interest) with GIST(Quattoni & Torralba, 2009), MM (Max-Margin)-background(Zhu et al., 2010), Object bank(Li et al., 2010), Improved Fisher Vector (IFV)(Perronnin et al., 2010), Laplacian Sparse coding SPM (LscSPM)(ShenghuaGao & Liang-TienChia, 2010), CENsus TRansform hISTogram (CENTRIST) (Wu & Rehg, 2011), Reconfigurable Bag of Words (RBoW)(Parizi et al., 2012), Bag of Parts (BoP)(Juneja et al., 2013), multi-channel CENTRIST (mCENTRIST) (Xiao et al., 2014), Important Spatial Pooling Region (ISPR)(Lin et al., 2014), Oriented Texture Curves (OTC) (Margolin et al., 2014), Scale Invariant Feature Transform (SIFT) (Zeglazi et al., 2016), and so on. The feature extraction algorithms under such traditional computer vision-based algorithms emphasize the low-level details of the images such as colors, intensity, gradients, orientations, etc. In other words, such algorithms are mostly local details oriented, and therefore, more suitable for specific types of images such as texture images. They are usually not ideal to represent the complex types of images such as scene images. Also, the feature size extracted by such algorithms is mostly higher than recent high-level features.

Furthermore, deep learning-based methods (Zhang et al., 2017; Gong et al., 2014; Guo & Lew, 2016; Guo et al., 2018; Tang et al., 2017; Kuzborskij et al., 2016; Zhou et al., 2016; He et al., 2016; Bai et al., 2019) are found to have noticeably better classification accuracies than existing traditional methods. Recent deep learning-based algorithms for scene representations are: CNN-MOP (Gong et al., 2014), CNN-sNBNL (Kuzborskij et al., 2016), VGG (Zhou et al., 2016), ResNet152 (He et al., 2016), EISR (Zhang et al., 2017), G-MS2F (Tang et al., 2017), SBoSP-fusion (Guo & Lew, 2016), BoSP-Pre_gp (Guo et al., 2018), CNN-LSTM (Bai et al., 2019), and so on. Gong et al. (2014) and Kuzborskij et al. (2016) employed Caffe model (Jia et al., 2014) to achieve features from the $FC$-layer for the scene images classification purpose. Gong et al. extracted multi-scale order-less features, which were obtained by extracting the global activation features ($FC$-layer) for each scale of the images and aggregated using the Vector of Locally Aggregated Descriptors (VLAD) pooling method. Similarly, Kuzborskij et al. also feed the output of $FC$-layers into the Naive Bayes Nearest Neighborhood classifiers. Zhou et al. (2016) released a new places related dataset to train the popular deep learning model such as VGG model (Simonyan & Zisserman, 2014). This leads to a promising classification accuracy of the images, especially scene images. He et al. (2016) proposed a novel architecture for deep learning, which followed the residual concepts and outperformed the previous off-the-shelf deep learning models such as the VGG model (Simonyan & Zisserman, 2014), GoogleNet model (Szegedy et al., 2015), etc. Zhang et al. (2017) sliced an image into multiple sub-images using random slicing and extracted deep features for each slice. Deep features of each slice were concatenated as combined deep features of the corresponding image. Finally, the fusion of such combined deep features with tag-based features yielded the final features representing an image for the classification. Tang et al. (2017) introduced a score-fusion technique to provide the probability-based deep features. They employed three intermediate classification layers of the GoogleNet model (Szegedy et al., 2015) for the score fusion, which improved classification performance remarkably. Guo et al. (Guo & Lew, 2016; Guo et al., 2018)

7

adopted the VGG-16 model to extract the deep features by developing the concept of the bag of surrogate parts (BoSP). Their method provides features with a fixed-size, which is lower than the most of the existing methods for the scene image representation despite the prominent classification accuracy. Recently, Bai et al. (2019) established a new deep learning model by incorporating Convolutional Neural Networks (CNNs) with Long Short Term Memory networks (LSTMs). They believe that the ordered slices of images as a sequence problem could be solved by the LSTMs model on top of CNNs model for the scene image classification. Their method, thus, offers prominent classification accuracy on scene images.

*2.2. Context-based image representation*

There have been some recent works (Zhang et al., 2017; Wang & Mao, 2019; Sitaula et al., 2019b) under context-based image representation methods. The main motivation of using such features is the use of human knowledge scattered in the form of annotation/description on the web, based on which people may be able to distinguish confusing complex scene images. For this, Zhang et al. (2017) extracted the description of the related images on the web to design the bag of words (BoW) features directly. Their method suffers from not only the occurrence of outliers but also the curse of higher feature size. To tackle this limitation, Wang & Mao (2019) devised the concept of filter bank using the category labels of ImageNet and Places to filter out the outliers to some degree. However, their method do not filter out the outliers accurately. The main reason of it not being able to filter out outliers is the dependence on pre-defined category labels only for the filter banks. This results in retrieval of more task-generic tags which are not suitable for scene images. To address such problem raised in the previous works, Sitaula et al. (2019b) developed a novel domain-specific filter bank that extracts the tag-based features by leveraging the semantic similarity of tags with scene image category labels. Their method not only generates rich tag-based features with the use of such filter banks, but also reduces the feature size of an image significantly.
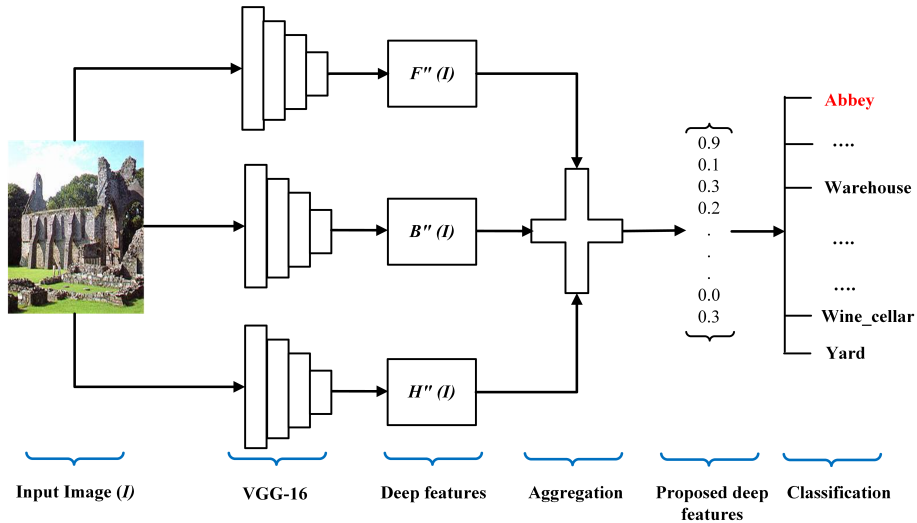
Figure 3: Overall pipeline of the proposed method. Three different pre-trained VGG-16 models yield foreground, background, and hybrid features represented by $F''(I)$, $B''(I)$ and $H''(I)$, respectively.

To sum up, the existing content-based and context-based image representation methods based on deep learning models outperform previous methods in most cases. This motivates us to explore further in content-based image representation based on the deep learning models, to achieve better representation of scene images. However, such methods in the literature have two main limitations. First, the existing methods consider either foreground or background information to represent the scene images, which may not be sufficient clue to some of the scene images requiring additional information such as hybrid information as a discriminating clue (see in Fig. 1). Second, scene images may contain three different types of information including foreground, background and hybrid information, which are complementary to each other in scene image representation tasks. Thus, we propose to fuse three different types of information (foreground, background and hybrid) that describe distinguishing clues in the scene representation. Our features bolster the classification performance significantly.
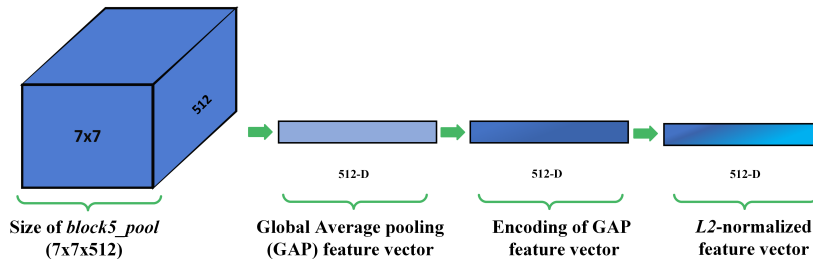
9

Figure 4: Three steps to achieve normalized feature vector after global average pooling (GAP) operation on the $5^{th}$ pooling layer ($p\_5$) followed by encoding step and normalization step.

## 3. Proposed method

To extract the proposed features, we follow four steps: foreground features extraction (Section 3.1), background features extraction (Section 3.2), hybrid features extraction (Section 3.3), and their aggregation (Section 3.4). The overall pipeline of our proposed method is shown in Fig. 3. For each of the first three steps, we have (1) global average pooling (GAP) features extraction that helps to capture both higher and lower activation values appropriate for scene representation, (2) encoding and (3) normalization. It is shown in Fig. 4. All the normalized feature vectors achieved are aggregated to achieve our proposed features as the final scene image representation.

### 3.1. Foreground features extraction

Foreground features often capture the object-based information in the scene images. These features are extracted from deep learning models pre-trained on ImageNet which consists of object images of $1,000$ categories. There are several pre-trained models such as ResNet-50 (He et al., 2016), GoogleNet (Szegedy et al., 2015), etc. for the foreground features extraction; nevertheless, we use VGG-16 model that has been frequently used in scene representation tasks (Guo & Lew, 2016; Guo et al., 2018; Bai et al., 2019) due to its simplicity and prominent features extraction capability. We represent VGG-16 model pre-trained on ImageNet as $F$. Here, Eq. (1) extracts GAP (Global Average Pooling) features

10

$F(I)$ from the $p\_5$ layer of the model $F$ (see details in Sec. 7). For the extraction of such GAP features, we average each feature map with height $(h)$ and width $(w)$. Similarly, the depth $(d)$ represents the total number of feature maps in the particular tensor. For instance, the $5^{th}$ pooling layer $(p\_5)$ of VGG-16 model has a three dimensional tensor of height $(h)$, width $(w)$, and depth $(d)$ as 7, 7 and 512, respectively. This results in the features size of 512-D after GAP operation on each feature map of the tensor. For this, we assume that the symbol $F_j^i$ represents the $i^{th}$ activation value of the $j^{th}$ feature map of the tensor.

$$F(I) = \frac{1}{h*w} * \left\{ \sum_{i=1}^{h*w} F_1^i, \sum_{i=1}^{h*w} F_2^i, , \cdots \sum_{i=1}^{h*w} F_d^i \right\}, \tag{1}$$

GAP features $(F(I))$ achieved from Eq. (1) are represented by the vector elements such as $\{f_1, f_2, \cdots f_d\}$, where $d$ is the size of such features and $f_j = \frac{\sum_{i=1}^{h*w} F_j^i}{h*w}$. These features are encoded as suggested by Guo et al. (2018) and Guo & Lew (2016), which has been found prominent for feature map encoding during the foreground based features extraction. Rather than utilizing such encoding for each feature map in Guo et al. (2018) and Guo & Lew (2016) , we employ it in our GAP features, which yields $F'(I)$ shown in Eq. (2). GAP features provide the discriminating information of scene images because it helps leverage both higher and lower activation values, which are discriminating clues of scene images classification.

$$F'(I) = \begin{cases} 0 & \text{if } f_j < mean(F(I)), \\ \frac{f_j}{max(F(I))} & \text{if } f_j \geq mean(F(I)). \end{cases}, \tag{2}$$

The encoded features are normalized using $L2$-norm to obtain $F''(I)$ as shown in Eq. (3). While doing normalization, we add epsilon (i.e., $\epsilon = 1e - 7$) with the denominator to avoid the divide-by-zero exception.

$$F''(I) = \frac{F'(I)}{\|F'(I)\|_2 + \epsilon}, \tag{3}$$

Eq. (3) yields the foreground features $(F''(I))$ of the input image.

*3.2. Background features extraction*

Background features represent the global layout information present in the images. These features are extracted from the deep learning model pre-trained on Places that involves background images of 365 categories. We represent VGG-16 model pre-trained on Places as $B$ in this work. The GAP features extracted from this model is shown in Eq. (4). Let $B_j^i$ represent the $i^{th}$ activation value for the corresponding $j^{th}$ feature map, and $B(I)$ represent the GAP features extracted from $B$.

$$B(I) = \frac{1}{h*w} * \left\{ \sum_{i=1}^{h*w} B_1^i, \sum_{i=1}^{h*w} B_2^i, , \cdots \sum_{i=1}^{h*w} B_d^i \right\}. \tag{4}$$

GAP features $B(I)$, which are represented by $\{b_1, b_2, \cdots, b_d\}$, are further encoded and normalized in similar ways to Eq. (2) and (3), respectively. Finally, the resulting features vector, say $B''(I)$, contains the background features.

*3.3. Hybrid features extraction*

Hybrid features represent the mixed features that are achieved from the deep learning model pre-trained on hybrid (ImageNet+Places) datasets. The datasets consist of combined images of objects and scenes under $1,365$ categories. For the extraction of such features, we also adopt the GAP features extracted from the $p\_5$ layer (Eq. (5)). Here, $H_j^i$ denotes the $i^{th}$ activation value of the corresponding $j^{th}$ feature map. $H(I)$ represents the GAP features extracted from $H$, where $H$ is the VGG-16 pre-trained model on the hybrid datasets.

$$H(I) = \frac{1}{h*w} * \left\{ \sum_{i=1}^{h*w} H_1^i, \sum_{i=1}^{h*w} H_2^i, , \cdots \sum_{i=1}^{h*w} H_d^i \right\}. \tag{5}$$

We again encode and normalize such a GAP features vector, which is represented by $\{h_1, h_2, \cdots, h_d\}$, extracted from Eq. (5) using Eqs. (2) and (3), respectively. This produces the hybrid features, say $H''(I)$.

---
**Algorithm 1** Extraction of the proposed features
---
**Input:** $F \leftarrow$ VGG-16 pre-trained on ImageNet database,

    $B \leftarrow$ VGG-16 pre-trained on Places database,

    $H \leftarrow$ VGG-16 pre-trained on Hybrid database,

    $I \leftarrow$ Input image for the feature extraction.

**Output:** $P(I) \leftarrow \{\}$ {Proposed features of image $I$}.

 1: Extract foreground features $(F''(I))$ of $I$ using Section 3.1.

 2: Extract background features $(B''(I))$ of $I$ using Section 3.2.

 3: Extract hybrid features $(H''(I))$ of $I$ using Section 3.3.

 4: Perform aggregation to achieve $P(I)$ using Section 3.4.

 5: **return** $P(I)$
---

*3.4. Aggregation*

These three types of deep features are aggregated to achieve our proposed features for the scene image representation. There are various simple yet efficient aggregation methods including $Min$, $Max$, $Mean$, and $Concat$. To this end, it is necessary to first find out the best and suitable aggregation method. For the selection of best aggregation method, we perform ablative analysis on various aggregation methods in Section 4.6. It shows that the feature size becomes 512-D for the Min, Max, and Mean methods, whereas it becomes $1,536$-D for the $Concat$ method. We choose to use the $Concat$ method on both datasets because it helps to represent three different types of information more accurately than other three methods and enable the state-of-the-art classification performance. Mathematically, the aggregation of the three different deep features including $F''(I)$, $B''(I)$, and $H''(I)$ is shown in Eq. (6).

$$P(I) = [F''(I), B''(I), H''(I)], \tag{6}$$

Alg. 1 lists the procedure to compute our features for the scene images representation.

## 4. Experiment and analysis

In this section, we discuss the experimental settings and compare our method with the previous methods, and perform ablative analysis of various parameters/components in the proposed method.

### 4.1. Datasets

For the experiments, we use two commonly used benchmark scene image datasets: MIT-67 (Quattoni & Torralba, 2009) and SUN-397 (Xiao et al., 2010). Both datasets contain numerous challenging images, involving within-class variations and between-class similarities which increase challenges for the classification performance.

**MIT-67** contains $15,620$ images under 67 categories. Each category includes at least 100 images. Some example images of this dataset are shown in Fig. 5. For the training and testing splits, we use the same split defined by Quattoni & Torralba (2009), which has been frequently used by the existing research (Quattoni & Torralba, 2009; Zhu et al., 2010; Li et al., 2010; Parizi et al., 2012; Juneja et al., 2013; Margolin et al., 2014; Lin et al., 2014; Zhang et al., 2017; Gong et al., 2014; Zhou et al., 2016; He et al., 2016; Guo & Lew, 2016; Guo et al., 2018; Tang et al., 2017; Bai et al., 2019; Kim, 2014; Wang & Mao, 2019; Sitaula et al., 2019b) . In particular, 80 images per category are used as training and 20 images per category are used as the testing in the experiments, which is a standard split defined by Quattoni & Torralba (2009).



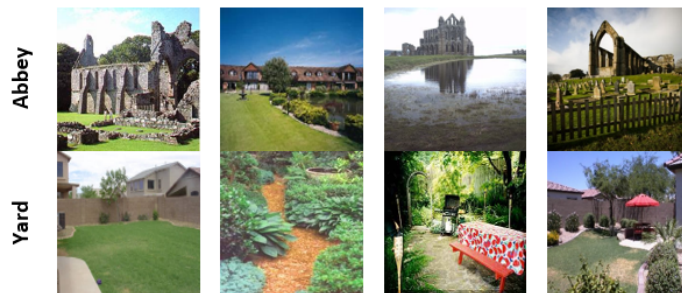Figure 5: Example images sampled from MIT-67 (Quattoni & Torralba, 2009).

Figure 6: Example images sampled from SUN-397 (Xiao et al., 2010).

**SUN-397** contains $108,754$ images under $397$ categories, where each category involves at least $100$ images. Some example images of this dataset are shown in Fig. 6. For training and testing, we use exactly the same splits defined by Xiao et al. (2010), which consists of 10 sets of train/test splits for the experiments. Several research works (Xiao et al., 2010; Margolin et al., 2014; Sánchez et al., 2013; Gong et al., 2014; Zhou et al., 2014; Simonyan & Zisserman, 2014; Yang & Ramanan, 2015; Wu et al., 2015; Dixit et al., 2015; Guo & Lew, 2016; Guo et al., 2018; Bai et al., 2019) have been tested on this dataset using content-based features extraction methods. To report the accuracy on this dataset, the mean accuracy of 10 splits is used, similar to previous research. In each split defined by Xiao et al. (2010), 50 images per category are used for training and 50 images per category are used for testing. The total number of sampled images used in the experiments for all 10 sets is $397,000$ ($10 \times 39,700$).

*4.2. Implementation*

To implement our method, we use Keras (Chollet et al., 2015) in Python (Python Core Team, 2015) language. Keras is used to implement pre-trained deep learning models (Simonyan & Zisserman, 2014; Kalliatakis, 2017)) to extract foreground, background and hybrid information of the scene images. For the classification purpose, we use the $L2$-Regularized Logistic Regression ($LR$) Classifier under the LibLinear (Fan et al., 2008) package in Python. To tune the best cost parameters ($C$) automatically, we perform grid search in the range

$\{1, 2, \cdots 50\}$ while keeping other parameters as default. All the experiments are conducted on a laptop with a NVIDIA GeForce GTX 1050 GPU.

*4.3. Comparison with state-of-the-art methods*

We report the accuracies in two different tables: Table 1 for MIT-67 and Table 2 for SUN-397.

Table 1 reports the accuracies of the previous methods and our method on the MIT-67 dataset. To minimize the bias, we only report the accuracies of such methods that are published using such datasets. We also report the method types (content-based and context-based) in the tables. For the methods belonging to content-based and context-based methods, we chose two methods (earliest one, and the latest one) for the comparison. We notice that our method produces 56.2% higher classification accuracy than the ROI with GIST (Quattoni & Torralba, 2009) and 1.8% higher than the most recent method (CNN-LSTM (Bai et al., 2019) method under the content-based features extraction methods). Furthermore, our method is 29.9% higher than the BoVW method and 5.9% higher than the TSF (Sitaula et al., 2019b) under the context-based features extraction methods. It shows that our method provides a significantly higher accuracy (**82.3**%) on this dataset outperforming both types of methods (content and context features extraction methods).

Table 2 presents the published accuracies of the previous methods as well as ours on the SUN-397 dataset. To date, there are not any context-based features extraction methods performed on this dataset. This may be because of the huge amount of images in this dataset that require heavy computation while performing query search on the web to achieve the context-based features. So we compare our method with some recent existing methods which belong to the content-based feature extraction methods. We observe that our method is 28% higher than Xiao et al. (2010), which is the very first method performed on this dataset, and 3.3% higher than the recent deep learning-based method, S-BoSP-Pre_gp (Guo et al., 2018). This apparently discloses the efficacy of our method which produces a noticeably higher accuracy (**66.3%**) on this huge benchmark

16

Table 1: Comparison of our proposed method using classification accuracy (%) with the previous methods on the MIT-67 dataset. Best results are bold face.

| Method | Accuracy(%) |
| --- | --- |
| Content-based feature extraction methods | |
| ROI with GIST (Quattoni & Torralba, 2009) | 26.1 |
| MM-background (Zhu et al., 2010) | 28.3 |
| Object Bank (Li et al., 2010) | 37.6 |
| RBoW (Parizi et al., 2012) | 37.9 |
| BOP (Juneja et al., 2013) | 46.1 |
| OTC (Margolin et al., 2014) | 47.3 |
| ISPR (Lin et al., 2014) | 50.1 |
| EISR (Zhang et al., 2017) | 66.2 |
| CNN-MOP (Gong et al., 2014) | 68.0 |
| VGG (Zhou et al., 2016) | 75.3 |
| ResNet152 (He et al., 2016) | 77.4 |
| S-BoSP-fusion (Guo & Lew, 2016) | 77.9 |
| BoSP-Pre_gp (Guo et al., 2018)) | 78.2 |
| G-MS2F (Tang et al., 2017) | 79.6 |
| CNN-LSTM (Bai et al., 2019) | 80.5 |
| Context-based feature extraction methods | |
| BoW (Wang & Mao, 2019) | 52.5 |
| CNN (Kim, 2014) | 52.0 |
| s-CNN(max) (Wang & Mao, 2019) | 54.6 |
| s-CNN(avg) (Wang & Mao, 2019) | 55.1 |
| s-CNNC(max) (Wang & Mao, 2019) | 55.9 |
| TSF (Sitaula et al., 2019b) | 76.5 |
| **Ours** | **82.3** |

Table 2: Comparison of our proposed method using classification accuracy (%) with the previous methods on the SUN-397 dataset.

| Method | Accuracy(%) |
|---|---|
| Content-based feature extraction methods | |
| Xiao et al. (Xiao et al., 2010) | 38.0 |
| OTC (Margolin et al., 2014) | 49.6 |
| FV (Sánchez et al., 2013) | 47.2 |
| CNN-MOP (Gong et al., 2014) | 51.9 |
| Places-CNN ((Zhou et al., 2014) | 54.3 |
| Hybrid-CNN (Zhou et al., 2014) | 53.8 |
| VGG-Net (Simonyan & Zisserman, 2014) | 51.9 |
| VGG-Net-DAG (Yang & Ramanan, 2015) | 56.2 |
| Metaforeground-CNN (Wu et al., 2015) | 58.1 |
| SFV-Places (Dixit et al., 2015) | 61.7 |
| S-BoSP-fusion (Guo & Lew, 2016) | 62.9 |
| S-BoSP-Pre_gp (Guo et al., 2018) | 63.7 |
| CNN-LSTM (Bai et al., 2019) | 63.0 |
| **Ours** | **66.3** |

dataset.

While seeing the classification accuracies of our method on both datasets (MIT-67 and SUN-397), we notice that our method produces competing and stable performance. Furthermore, regarding the feature size, our method has a lower dimensional size, which is just $1,536$-D; however, the main contender of our method (S-BoSP-Pre_gp) on the SUN-397 dataset has $9,216$-D features size, which is over 6 times larger than ours. Similarly, the feature size of CNN-LSTM which is the main contender on the MIT-67 dataset, still has a greater feature size ($4,096$-D) than ours. In general, the classification time will increase with a growing size of features (Sitaula et al., 2019a). Therefore, our method consumes

a lower classification time than such contender methods.

### 4.4. Ablative study of pooling methods

The selection of appropriate pooling layer is essential in features extraction while using pre-trained deep learning models. We perform extensive experiments on both datasets using three pre-trained VGG-16 models. Specifically, we analyze five pooling layers of each of the VGG-16 models that has been trained on ImageNet, Places, and mixed data (ImageNet+Places). The five pooling layers used are $block1\_pool$, $block2\_pool$, $block3\_pool$, $block4\_pool$ and $block\_pool$, which are represented by $p\_1$, $p\_2$, $p\_3$, $p\_4$, and $p\_5$, respectively. The experimental results performed on both datasets are shown in Fig. 7. In the figure, we notice that the appropriate pooling layer for the distinguishing features extraction is the $5^{th}$ pooling layer ($p\_5$). The classification accuracy using such a layer on both datasets is higher than other layers. Thus, we deduce that the $5^{th}$ pooling layer extracts the high-level information of the image that has the ability for better representation of scene images. This leads us to utilizing this layer to achieve the corresponding information including foreground, background and hybrid information of the scene images.

### 4.5. Ablative study of individual features

We study the contribution of individual features used in our method on both datasets to see how they affect classification individually. Additionally, this study helps us to learn the highly promising features type among three different features in the scene image representation. The individual features include foreground features ($F''(I)$), background features ($B''(I)$) and hybrid features ($H''(I)$). The classification accuracies using three individual features achieved on both datasets are illustrated in Fig. 8. While observing the bar graph, features based on background information (background features) are better than the remaining two different types of features for the classification. Also, the superior accuracy of hybrid features than foreground features further unveil that hybrid features are also more important than foreground features on both
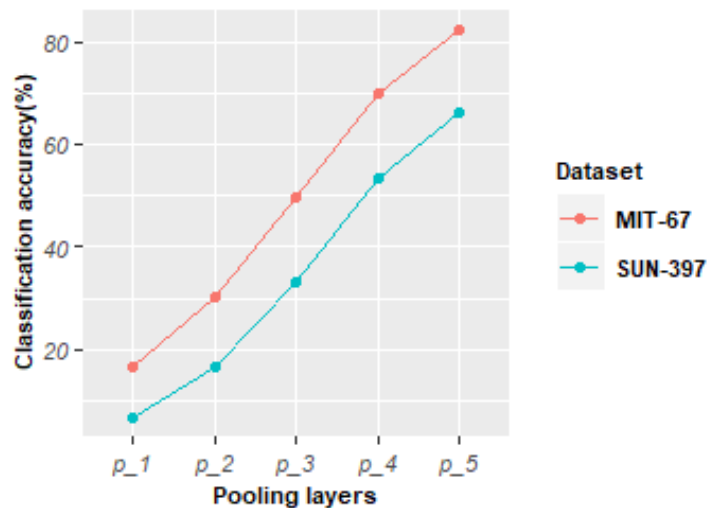
Figure 7: Comparative analysis of five pooling layers starting from $p\_1$ to $p\_5$ on MIT-67 and SUN-397 datasets.

datasets. Thus, we believe that the majority of the separability capability is attributed to the background information in most cases.

### 4.6. Ablative study of aggregation methods

Features aggregation is also one of the important steps in our method. For this, we perform experiments using four different aggregation methods including Minimum ($Min$), Maximum ($Max$), $Mean$ and Concatenation ($Concat$) on both datasets. The experimental result are shown in Fig. 9. Results reveal that the $Concat$ method outperforms all other methods on both datasets. We speculate that the $Concat$ method on uniform sized features alleviates bias during features aggregation, thereby preserving multi-information uniformly for the classification purpose. As such, we adopt this method in our work.

### 4.7. Ablative study of combined features

In this subsection, we analyze the efficacy of the combined features on both datasets using classification accuracy. For this, we combine three types of infor-
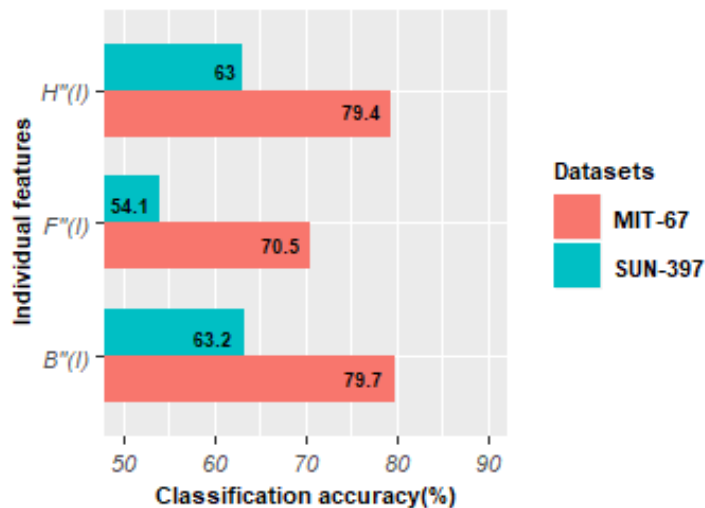
Figure 8: Comparative analysis of each individual features on the MIT-67 and SUN-397 datasets.

mation ($F''(I)$, $B''(I)$, and $H''(I)$) and provides four total combinations including $[F''(I), B''(I)]$, $[F''(I), H''(I)]$, $[B''(I), H''(I)]$, and $[F''(I), B''(I), H''(I)]$. For the combination of features, we use the *Concat* aggregation method as suggested in Sec. 9 and present the results in Table 3. We see that the fusion of three types of features outperforms all other combinations on both datasets, in terms accuracy. We conjecture that these three types of features are complementary so that the fusion empowers a better representation of scene images.

*4.8. Computation time*

We analyze the computation time (seconds) and list the results in Table 4. For SUN-397, we average the computation time of all 10 sets. We observe that the average features extraction time per image for all the images including training and testing sets on both SUN-397 (39,7000 images) and MIT-67 (6,700 images) is 0.1 seconds. Similarly, the average classification time of the testing images on the SUN-397 (19,850 images) and MIT-67 (1,340 images) are 0.0004 seconds and 0.0005 seconds, respectively.
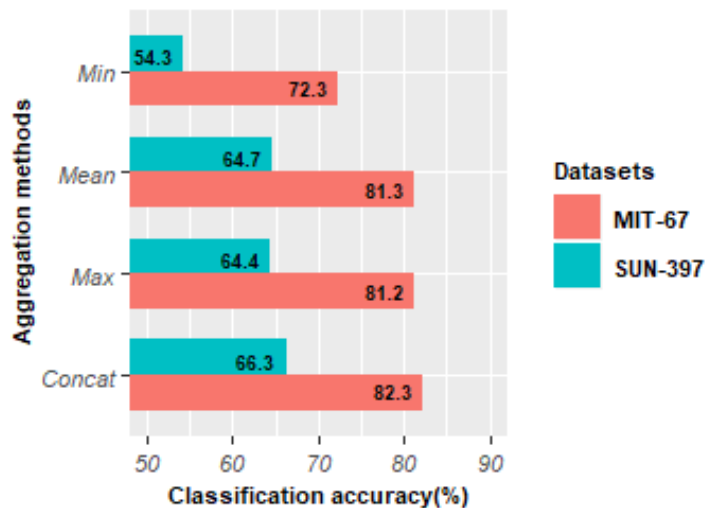
Figure 9: Comparative analysis of four aggregation methods on the MIT-67 and SUN-397 datasets.

## 5. Conclusion and future works

In this paper, we have proposed a method that aggregates three different types of deep features for the scene image representation. Experimental results on the commonly-used benchmark scene datasets demonstrate a better classification performance of our method than the state-of-the-art methods. Furthermore, our method also outputs a noticeably lower size of features of the scene images.

Our proposed method is more suitable for scene images than other types of images, because all our captured information is focused on scene images rather than images such as satellite images, biomedical images, Internet of Things (IoT) images. Those images may require other discriminating clues such as texture, global layout, temporal information, spatial information, and so on for representation. In the future, we would like to investigate other types of images for better representations.

Table 3: Comparative analysis (%) of the combined features on two datasets.

| Comb. layers | MIT-67 (%) | SUN-397 (%) |
|---|---|---|
| $[F''(I), B''(I)]$ | 81.7 | 65.1 |
| $[F''(I), H''(I)]$ | 80.2 | 63.6 |
| $[B''(I), H''(I)]$ | 81.5 | 65.8 |
| $[F''(I), B''(I), H''(I)]$ | **82.3** | **66.3** |

Table 4: Computation time (seconds) consumed by three main steps such as feature extraction step, training step, and testing step for our proposed method on two datasets.

| Dataset | Feat. extraction step | Training step | Testing step |
|---|---|---|---|
| MIT-67 | 756.7 | 7.8 | 0.8 |
| SUN-397 | 4779.4 | 67.8 | 8.1 |

## References

Bai, S., Tang, H., & An, S. (2019). Coordinate cnns and lstms to categorize scene images with multi-views and multi-levels of abstraction. *Expert Syst. Appl.*, *120*, 298–309.

Chollet, F. et al. (2015). Keras. `https://github.com/fchollet/keras`.

Dalal, N., & Triggs, B. (2005). Histograms of oriented gradients for human detection. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)* (pp. 886–893).

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). ImageNet: a large-scale hierarchical image database. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*.

Dixit, M., Chen, S., Gao, D., Rasiwasia, N., & Vasconcelos, N. (2015). Scene classification with semantic fisher vectors. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)* (pp. 2974–2983).

Fan, R.-E., Chang, K.-W., Hsieh, C.-J., Wang, X.-R., & Lin, C.-J. (2008). Lib-linear: A library for large linear classification. *Journal of Machine Learning Research*, *9*, 1871–1874.

Gong, Y., Wang, L., Guo, R., & Lazebnik, S. (2014). Multi-scale orderless pooling of deep convolutional activation features. In *Proc. Eur. Conf. Comput. Vis. (ECCV)* (pp. 392–407).

Guo, Y., & Lew, M. S. (2016). Bag of surrogate parts: one inherent feature of deep cnns. In *Proc. BMVC*.

Guo, Y., Liu, Y., Lao, S., Bakker, E. M., Bai, L., & Lew, M. S. (2018). Bag of surrogate parts feature for visual recognition. *IEEE Trans. Multimedia*, *20*, 1525–1536.

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (pp. 770–778).

Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadar-rama, S., & Darrell, T. (2014). Caffe: Convolutional architecture for fast feature embedding. In *Proc. 22nd ACM Int. Conf. Multimedia* (pp. 675–678).

Juneja, M., Vedaldi, A., Jawahar, C., & Zisserman, A. (2013). Blocks that shout: Distinctive parts for scene classification. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)* (pp. 923–930).

Kalliatakis, G. (2017). Keras-vgg16-places365. `https://github.com/GKalliatakis/Keras-VGG16-places365`.

Kim, Y. (2014). Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, .

Kuzborskij, I., Maria Carlucci, F., & Caputo, B. (2016). When naive bayes nearest neighbors meet convolutional neural networks. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (pp. 2100–2109).

Lazebnik, S., Schmid, C., & Ponce, J. (2006). Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* (pp. 2169–2178).

Li, L.-J., Su, H., Fei-Fei, L., & Xing, E. P. (2010). Object bank: A high-level image representation for scene classification & semantic feature sparsification. In *Proc. Adv. Neural Inf. Process. Syst. (NIPS)* (pp. 1378–1386).

Lin, D., Lu, C., Liao, R., & Jia, J. (2014). Learning important spatial pooling regions for scene classification. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)* (pp. 3726–3733).

Margolin, R., Zelnik-Manor, L., & Tal, A. (2014). Otc: A novel local descriptor for scene classification. In *Proc. Eur. Conf. Comput. Vis. (ECCV)* (pp. 377–391).

Oliva, A. (2005). Gist of the scene. In *Neurobiology of Attention* (pp. 251–256). Elsevier.

Oliva, A., & Torralba, A. (2001). Modeling the shape of the scene: a holistic representation of the spatial envelope. *Int. J. Comput. Vis.*, *42*, 145–175.

Parizi, N., Oberlin, J. G., & Felzenszwalb, P. F. (2012). Reconfigurable models for scene recognition. In *Proc. Comput. Vis. Pattern Recognit.(CVPR)* (pp. 2775–2782).

Perronnin, F., Sánchez, J., & Mensink, T. (2010). Improving the fisher kernel for large-scale image classification. In *Proc. European Conference on Computer vision (ECCV)* (pp. 143–156).

Python Core Team (2015). *Python: A Dynamic, Open Source Programming language*. Python Software Foundation.

Quattoni, A., & Torralba, A. (2009). Recognizing indoor scenes. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)* (pp. 413–420).

Sánchez, J., Perronnin, F., Mensink, T., & Verbeek, J. (2013). Image classification with the fisher vector: Theory and practice. *International Journal of Computer Vision (IJCV)*, *105*, 222–245.

ShenghuaGao, I.-H., & Liang-TienChia, P. (2010). Local features are not lonely–laplacian sparse coding for image classification, . (pp. 3555–3561).

Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, .

Sitaula, C., Xiang, Y., Aryal, S., & Lu, X. (2019a). Unsupervised deep features for privacy image classification. In *Proc. Pacific-Rim Symposium on Image and Video Technology (PSIVT)* (pp. 404–415).

Sitaula, C., Xiang, Y., Basnet, A., Aryal, S., & Lu, X. (2019b). Tag-based semantic features for scene image classification. In *Proc. Int. Conf. on Neural Inf. Process. (ICONIP)* (pp. 90–102).

Sitaula, C., Xiang, Y., Zhang, Y., Lu, X., & Aryal, S. (2019c). Indoor image representation by high-level semantic features. *IEEE Access*, *7*, 84967–84979.

Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., & Rabinovich, A. (2015). Going deeper with convolutions. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)* (pp. 1–9).

Tang, P., Wang, H., & Kwong, S. (2017). G-MS2F: GoogLeNet based multistage feature fusion of deep CNN for scene recognition. *Neurocomputing*, *225*, 188 – 197.

Wang, D., & Mao, K. (2019). Task-generic semantic convolutional neural network for web text-aided image classification. *Neurocomputing*, *329*, 103–115.

Wu, J., & Rehg, J. M. (2011). CENTRIST: A visual descriptor for scene categorization. *IEEE Trans. Pattern Anal. Mach. Intell.*, *33*, 1489–1501.

Wu, R., Wang, B., Wang, W., & Yu, Y. (2015). Harvesting discriminative meta objects with deep cnn features for scene classification. In *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)* (pp. 1287–1295).

Xiao, J., Hays, J., Ehinger, K. A., Oliva, A., & Torralba, A. (2010). Sun database: Large-scale scene recognition from abbey to zoo. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)* (pp. 3485–3492).

Xiao, Y., Wu, J., & Yuan, J. (2014). mCENTRIST: a multi-channel feature generation mechanism for scene categorization. *IEEE Trans. Image Process.*, *23*, 823–836.

Yang, S., & Ramanan, D. (2015). Multi-scale recognition with dag-cnns. In *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)* (pp. 1215–1223).

Zeglazi, O., Amine, A., & Rziza, M. (2016). Sift descriptors modeling and application in texture image classification. In *Proc. 13th Int. Conf. Comput. Graphics, Imaging and Visualization (CGiV)* (pp. 265–268).

Zhang, C., Zhu, G., Huang, Q., & Tian, Q. (2017). Image classification by search with explicitly and implicitly semantic representations. *Information Sciences*, *376*, 125–135.

Zhou, B., Khosla, A., Lapedriza, A., Torralba, A., & Oliva, A. (2016). Places: An image database for deep scene understanding. *arXiv preprint arXiv:1610.02055*, .

Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., & Torralba, A. (2017). Places: A 10 million image database for scene recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, *40*, 1452–1464.

Zhou, B., Lapedriza, A., Xiao, J., Torralba, A., & Oliva, A. (2014). Learning deep features for scene recognition using places database. In *Proc. Adv. Neural Inf. Process. Syst. (NIPS)* (pp. 487–495).

Zhu, J., Li, L.-j., Fei-Fei, L., & Xing, E. P. (2010). Large margin learning of upstream scene understanding models. In *Proc. Adv. Neural Inf. Process. Syst. (NIPS)* (pp. 2586–2594).