

# Deep Scene Image Classification with the MFAFVNe

Yunsheng Li Mandar Dixit Nuno Vasconcelos University of California, San Diego La Jolla, CA 92093

yul554@ucsd.edu mdixit@ucsd.edu nvasconcelos@ucsd.edu

work trained for object recognition to the task of scene image classification is considered. An embedded implementation of the recently proposed mixture of factor analyzers tion must be invariant to the locations of individual objects *Fisher vector (MFA-FV) is proposed. This enables the de-* which can change drastically within the same scene class sign of a network architecture, the MFAFVNet, that can be Eventually, sophisticated pooling strategie, such as the vectrained in an end to end manner. The new architecture in-tor of locally aggregated descriptors (VLAD) [12] or the volves the design of a MFA-FV layer that implements a sta-Fisher vector (FV) [23] emerged as the dominant pooling *tistically correct version of the MFA-FV, through a combi-* mechanisms for scene classification. nation of network computations and regularization. When compared to previous neural implementations of Fisher vec-tors of choice for scene classification. In fact, the impletors, the MFAFVNet relies on a more powerful statistical mentation of a holistic scene classifier with a CNN is almodel and a more accurate implementation. When compared to previous non-embedded models, the MFAFVNet images. The main challenges are the assembly of a large relies on a state of the art model, which is now embedded dataset of such images and the standard difficulties of traininto a CNN. This enables end to end training. Experiments ing a deep network. These problems have been addressed show that the MFAFVNet has state of the art performance in [29], through the assembly of the Places dataset and its on scene classification.

## 1. Introduction

applications of computer vision such as robotics, image ferred across datasets, by simple finetuning, is one of their search, geo-localization, etc. It is also a challenging problem, e.g. object recognition methods do not necessarily work for scenes. This is because scenes include both a ognize objects) to a holistic task (classification of whole holistic component, the *gist* of the scene, and an object-scenes) cannot, in general, be solved by simple finetuning. based component. Furthermore, the object vocabulary is The previous success of pooling on this type of trans usually open-ended and it does not suffice to recognize objects, as most scenes are collections of objects in charac-deep learning realm. Several authors proposed Fisher vecteristic spatial layouts. There is also a need to model relationships between objects. Historically, this motivated different approaches to scene classification, including holistic like approach, based on the extraction of features from in-

tual relationships between objects [14] or semantics prop-ing [17] and finally used to implement descriptors such as erties such as the objects in the scene [15] or scene at-the VLAD [11] or Fisher vector [17]. Later, [5] proposed tributes [25]. Many of these approaches were based on the the semantic Fisher Vector, which extracts features from the

bag-of-features (BoF) representation, using local features such as SIFT or HOG [19, 3], combined through a pool-The problem of transferring a deep convolutional net-ing operator. Pooling has a critical role in scene classifi-

In recent years, CNNs have become the feature extracuse to train a deep scene classifier. The use of deep *object*based representations for scene classification has, however proven more challenging. This is, in great part, because object-based scene classification is a difficult transfer learn-Scene image classification is an important problems for ing problem. While the ease with which CNNs can be transgreatest assets for vision, this procedure has its limitations The transfer from a localized representation (needed to rec-

gist descriptors [20] and descriptors based on local features. termediate CNN layers, which were then fed to dictionary Localized approaches included descriptors of contex- learning methods such as clustering [11] or sparse cod-

softmax layer at the top of the CNN, converting features proximations can be quite sub-optimal. from probability space to the natural parameter space. An alternative strategy, proposed by [6], is to use a better

however, prevents end to end training and, consequently, discussed above. the finetuning of the object network to the scene classifica- In this work, we derive an embedded implementation

Since CNN features are high dimensional, it is impractical to rely on Gaussians of full covariance. Instead, the mixture components are chosen to have diagonal covariance. This creates problems when the feature manifold is nonponents are required and the Fisher vector is very high dimensional. This is indeed the norm for computer vision applications, where Fisher vectors usually have several thousand dimensions. While the CNN is trained to produce linearly separable responses to the different classes, there is no guarantee CNN feature distributions are prone to modon the diagonal-GMM is likely to be very high dimensional classification. and potentially sub-optimal for scene classification.

Recently, some works have attempted to solve these problems. One possibility is to bypass the Fisher vector In this section, we review the main ideas behind the altogether. For example, [7] proposed a compact bilinear MFA-FV. pooling (CBP) mechanism that enables end-to-end training by simple backpropagation. While this was shown applica- **2.1. Mixture of Factor Analyzers** ble to scene classification, the performance of CBP is inferior to those of previous approaches, such as the semantic Fisher vector of [5] or the sparse coding methods of [18], for equivalent object CNNs. Another possibility is to embed the Fisher vector in the CNN architecture, by deriving a neural can be explained by a small number d < D of hidden or lanetwork implementation of its equations. [1] proposed the tent *factors*, usually represented as a factor vector  $z \in \mathbb{R}^d$ . NetVLAD, an embedded implementation of the VLAD de-Observations x are assumed to be sampled according to the scriptor, and [26] proposed the Deep FisherNet, an embedded implementation of the GMM Fisher vector. However, to avoid the difficulties of the complete Fisher vector, these methods make approximations, such as disregarding covari-where  $\mu$  is the mean value of  $x, \Lambda \in \mathbb{R}^{D \times d}$  is a fac-

All these methods suffer from two drawbacks. First, the model of CNN feature statistics than the diagonal GMM. Fisher vector structure is not easy to integrate in a CNN. Instead, this work proposed a Fisher vector based on the This is because the Fisher vector is defined with respect mixture of factor analyzers (MFA). This has the advantage to the probability distribution of the CNN features, usually of accounting for the covariance information of each mixestimated with a mixture learned by maximum likelihood. ture component, which is modeled through factor analysis. The Fisher vector is then a complex expression of the mix-Under the MFA model, good results can be achieved with ture parameters, which changes when these change. In result, the Fisher vector cannot be learned by simply back- dimension. However, while the MFA-Fisher vector (MFApropagating the output of the scene classifier. All methods FV) currently holds state of art results for scene classificaabove avoid this difficulty by using the CNN to extract fea- tion, it is not an integrated model, i.e. it is learned indetures and learning the Fisher vector indepeddently. This, pendently of the network. This raises all the reservations

tion task. In result, the transfer between the two tasks relies of the MFA-FV, and use it to design a network architecsolely on the Fisher vector, which is sub-optimal. ture, the MFAFVNet, which can be trained in an end to The second problem is that the Fisher vector is usually end manner. This involves the derivation of a MFA-FV learned with respect to a Gaussian mixture model (GMM). layer that implements a statistically correct version of the MFA-FV, through a combination of network computations and regularization. The computations replicate those of the MFA-FV, regularization guarantees that all parameters have a statistically valid interpretation. When compared to prelinear. In this case, a large number of diagonal GMM comon a more powerful statistical model, which accounts for covariance information, and a more accurate implementation. This results in significant performance gains for scene classification. When compared to previous non-embedded models, the MFAFVNet relies on a state of the art model, which is embedded. This enables end to end training and eling with the diagonal GMM. On the contrary, given the better scene classification performance. Extensive experihighly non-linear feature transformation implemented by a ments on the MIT Indoor and SUN datasets show that the deep CNN, this is unlikely. Hence, a Fisher vector based MFAFVNet achieves state of the art performance for scene

# 2. The MFA Fisher Vector

The factor analysis (FA) model is a probabilistic extension of principal component analysis (PCA) [23]. Given a vector  $x \in \mathbb{R}^D$  of D observations, it explains its covariance structure by assuming that the variability of the observations

$$x - \mu = \Lambda z + \epsilon, \tag{1}$$

ance structure (VLAD) or using crude approximations of tor loading matrix, and  $\epsilon$  is a noise term. Factors z and posterior mixture probabilities (Deep FisherNet). These apnoise  $\epsilon$  are independent of each other and Gaussian, and the



non-linear manfolds (right) can require many mixture components to approximate, when covariances are diagonal (shown in red). However, they ture model. In general, many components are needed to can frequently be approximated by a few MFA components (in green). achieve a good approximation of the distribution p(x). As

as  $\mathcal{N}(\epsilon; 0, I)$ . The factors can be dependent, i.e. they are dimensional subspaces. The MFA is a substantially better Gaussian with covariance

$$\Sigma = cov(x - \mu) = cov(\Lambda z + \epsilon) = \Lambda \Lambda^T + \psi. \quad (2)$$

Hence, the factor loading matrix  $\Lambda$  has a role similar to the principal component matrix of PCA.

The mixture of factor analyzers (MFA) is a mixture model whose components follow the factor analysis model. A MFA of C components is defined by the distributions

$$p(c) = \pi_c$$

$$p(z|c) = \mathcal{N}(z; 0, I)$$

$$p(x|z,c) = \mathcal{N}(x; \Lambda_c z + \mu_c, \Psi_c)$$

ponent is a FA of mean  $\mu_c$ , factor loading matrix  $\Lambda_c$  and is the state of the art representation for scene classification. noise covariance  $\Psi_c$ .

these model parameters so as to maximize the likelihoods this integration. of a set of observations  $\{x_i\}_{i=1}^N$ .

### 2.2. Fisher vectors

Given a dataset  $\mathcal{D} = \{x_i\}$  and a probability model  $p(x;\theta)$  the score  $\mathcal{G}(\theta) = \frac{\partial}{\partial \theta} \log p(\mathcal{D};\theta)$  measures the sensitivity of the likelihood  $p(\mathcal{D}; \theta)$  to parameter  $\theta$ . The nor- **3.1. The MFA-FV laye** malization of this gradient vector by the square root of the Fisher information matrix  $\mathcal{I}(\theta) = -\sum_i \frac{\partial^2}{\partial \theta^2} \log p(x_i; \theta)$ , network, we start by defining i.e. the vector  $\mathcal{I}^{-1/2}\mathcal{G}(\theta)$  is usually denoted as the *Fisher vector* [23]. However, because the Fisher information can be difficult to compute, it is frequently ignored and the

Fisher vector reduces to the score  $\mathcal{G}(\theta)$ . As is common in computer vision, we adopt this practice in this work. The Fisher vector is commonly used with the standard Gaussian mixture model, defined by

$$p(c) = \pi_c$$
(6)  
$$p(x|c) = \mathcal{N}(x; \mu_c, \Sigma_c).$$
(7)

However, because vision data tends to be high-dimensional, it is usually difficult to use full covariance Gaussians in this model, and the covariances  $\Sigma_c$  are assumed as diago-Figure 1. Probability distributions defined on linear susbpaces (left) or nal. This removes a lot of the expressiveness of the mixillustrated in Figure 1, this is particularly true when the data lives on correlated low-dimensional subspaces or nonnoise variables are assumed uncorrelated, i.e. distributed linear manifolds that can be approximated by a set of lowdistributed as  $\mathcal{N}(z; 0, \psi)$ , but the matrix  $\psi$  is sometimes as-model for this type of data since, in this case, only a few sumed diagonal. It follows from the linearity of (1) that x is mixture components and a small number d of hidden factors are required to estimate the covariance structure of (2). This is illustrated in Figure 1 as well. This observation, motivated the introduction of the MFA-Fisher vector (MFA-FV) in [6], which was shown to have the form

$$\mathcal{G}_{\mu_c}(\mathcal{I}) = \sum_i p(c|x_i;\theta)\psi^{-1}(I - \Lambda_c \Gamma_c)(x_i - \mu_c)(8)$$
  

$$\mathcal{G}_{\Lambda_c}(\mathcal{I}) = \sum_i p(c|x_i;\theta)\psi^{-1}(\Lambda_c \Gamma_c - I)$$
  

$$[(x_i - \mu_c)(x_i - \mu_c)^T \Gamma_c^T - \Lambda_c] \qquad (9)$$
  

$$\Gamma_c = \Lambda_c^T S_c^{-1}. \qquad (10)$$

<sup>(4)</sup> This work has also shown that the MFA-FV is a substan-(5) tially better representation than the classical FV when the observations x are feature vectors produced by a deep conwhere p(c) is the probability of component c and this com-volutional neural network (CNN). As far as we know, this However, [6] did not integrate the MFA-FV in the network The MFA can be learned with a EM algorithm, which is computation. This prevents end-to-end training and the tundiscussed in [8]. This iterates between an expectation step ing of the network to the scene classification task. Since end that computes expected values for the hidden class c and to end training is an important reason for the recent success factors z variables, and a maximization step that updates of the deep CNN architecture, it appears natural to pursue

# 3. Network implementation of the MFA-FV

In this section, we derive a neural network implementation of the MFA-FV.

To derive a version of (8)-(10) implementable as a neural

$$\Delta_{ic} = x_i - \mu_c \tag{11}$$
$$S_c = \Lambda_c \Lambda_c^T + \psi_c. \tag{12}$$



Figure 2. The MFA-FV layer. The expressions in red are the parameters of the  $c^{th}$  MFA component. The remaining expressions show what is computed at each stage of the network.

Combining this with (10),

$$\begin{split} \psi_c^{-1}(I - \Lambda_c \Gamma_c) &= (S_c - \Lambda_c \Lambda_c^T)^{-1} (I - \Lambda_c \Lambda_c^T S_c^{-1}) \\ &= S_c^{-1} (I - \Lambda_c \Lambda_c^T S_c^{-1})^{-1} (I - \Lambda_c \Lambda_c^T S_c^{-1}) \\ &= S_c^{-1}, \end{split}$$

and it follows that (8)-(9) can be written as

$$\mathcal{G}_{\mu_c}(\mathcal{I}) = \sum_i p(c|x_i;\theta) S_c^{-1} \Delta_{ic}$$
(13)
$$\mathcal{G}_{\Lambda_c}(\mathcal{I}) = -\sum_i p(c|x_i;\theta) S_c^{-1} [\Delta_{ic} \Delta_{ic}^T S_c^{-1} \Lambda_c - \Lambda_c]$$

$$= -\sum_i p(c|x_i;\theta) S_c^{-1} \Delta_{ic} [S_c^{-1} \Delta_{ic}]^T \Lambda_c$$

$$+ \sum_i p(c|x_i;\theta) S_c^{-1} \Lambda_c$$
(14)

Furthermore, since (2) implies that the  $c^{th}$  mixture component p(x|c) is distributed as  $\mathcal{N}(x, \mu_c, S_c)$ , it follows that

$$\begin{aligned}
\rho(c|x_i;\theta) &= \frac{\pi_c \mathcal{N}(x_i;\mu_c,S_c)}{\sum_k \pi_k \mathcal{N}(x_i;\mu_k,S_k)} \\
&= \frac{\frac{\pi_c}{|S_c|^{\frac{1}{2}}} \exp\{-\frac{1}{2}\Delta_{ic}^T S_c^{-1} \Delta_{ic}\}}{\sum_k \frac{\pi_k}{|S_k|^{\frac{1}{2}}} \exp\{-\frac{1}{2}\Delta_{ic}^T S_k^{-1} \Delta_{ic}\}}
\end{aligned}$$
(15)

Denoting

$$P_c = S_c^{-1},$$
  

$$\Omega_c = S_c^{-1}\Lambda_c,$$
  

$$\kappa_c = \frac{\pi_c}{|S|^{\frac{1}{2}}},$$

finally leads to

$$\mathcal{G}_{\mu_c}(\mathcal{I}) = \sum_i p(c|x_i; \theta) P_c \Delta_{ic}$$

$$p(c|x_i;\theta) = -\sum_{i} p(c|x_i;\theta) \{ P_c \Delta_{ic} (P_c \Delta_{ic})^T \Lambda_c - \Omega_c \}$$

$$p(c|x_i;\theta) = \frac{\kappa_c \exp\{-\frac{1}{2} \Delta_{ic}^T P_c \Delta_{ic}\}}{\sum_k \kappa_k \exp\{-\frac{1}{2} \Delta_{ik}^T P_k \Delta_{ik}\}}$$
(20)
(21)

An implementation of (20) as a neural network layer is shown in Figure 2. The bottom branch computes the posterior probability of (21). The top branch computes the remainder of the argument of the summation in (20). The computations of (19) are similar. The bottom branch is identical, and the top branch omits the operations beyond  $P_c \Delta_{ic}$ . However, preliminary experiments showed no gains for the addition of this component. Hence, we use only the second order information, i.e. (20). Note that the operations inside circle are applied entry-wise, the boxes implement matrix multiplications implementable with standard layers of weights, the outer product layer is similar to that of [16], and the dot-product layer can be implemented with an elementwise multiplication and a sum.

### **3.2. Relation to other Fisher vectors**

5) The MFA-FV is related to various previous representations of the same type. For example, if  $\Lambda_c$  is the identity matrix and  $S_c$  a diagonal matrix of elements  $\sigma_{ab}^2$ , then (19) reduces to

$$\mathcal{G}_{\mu_{ck}}(\mathcal{I}) = \sum_{i} p(c|x_{ik};\theta) \frac{(x_{ik} - \mu_{ck})}{\sigma_{ck}^2}$$
(22)

and (20) to

$$\mathcal{G}_{\sigma_{ck}}(\mathcal{I}) = \sum_{i} p(c|x_i;\theta) \frac{1}{\sigma_{ck}^2} \left\{ \frac{(x_{ik} - \mu_{ck})^2}{\sigma_{ck}^2} - 1 \right\}, \quad (23)$$

which are similar to the Fisher Score of the standard Gaussian mixture model [23, 21]. Further omitting the second other information leads to

$$\mathcal{G}_{\mu_{ck}}(\mathcal{I}) = \sum_{i} p(c|x_i;\theta)(x_{ik} - \mu_{ck})$$
(24)



Figure 3. Architecture of MFAFVNet.

and replacing the posterior probabilities  $p(c|x_i; \theta)$  by binary **3.4. Loss Function** variables  $a_{ic}$  such that  $a_{ic} = 1$  if  $\mu_c$  is the closest mean to  $x_i$  and  $a_{ic} = 0$  reduces to the VLAD

$$V_{ck} = \sum_{i} a_{ic} (x_{ik} - \mu_{ck}).$$
 (2)

Note that the variables  $a_{ic}$  can be obtained by replacing the softmax by a max in Figure 2.

### **3.3. The MFAFVnet**

Figure 3. A model pretrained on ImageNet is first used to favored. A similar observation can be made for (12), which extract a vector p(x) of image features. Our implementation is the only equation to constrain  $\psi_c$ . uses either Alexnet or VGG, both of which include a se- On the other hand, some of the relationships must be quence of several convolutional layers and three fully con- enforced to maintain the MFA-FV interpretation. These nected layers. As is common for scene classification, this are (16), (17), and the fact that  $P_c$  is a symmetric matrix. network is applied to image patches [23], producing multi- These relationships can be enforced by adding regularizaple feature maps per image to classify. When these patches tion terms to the loss function used to train the network. are of a single scale, the model is converted to a fully convo-Given a training set  $\mathcal{D} = \{(x_i, y_i)\}$  this leads to a loss funclutional network. When patches of multiple scales are used, tion the final pooling layer is replaced with a region-of-interest (ROI) pooling layer, which accepts feature maps of multiple sizes and produces a fixed size output to the fully connected layers. This is similar to the standard practice in the object detection literature [10, 9]. As is usual in the Fisher vector literature [23] the feature vector p(x) is subject to dimensionality reduction. This is implemented by a fully in Figure 2. Note that the MFA-FV layer pools multiple local features, corresponding to objects of different sizes and in different locations of the input image. It produces a single feature vector that represents the whole input image. As is standard in the Fisher vector literature [21], two normalization layers (power normalization and L2 normalization) where s(x) is the input to the softmax at the top of the

he parameters  $\mu_c$ ,  $P_c$ ,  $\Lambda_c$ ,  $\Omega_c$ , and  $\log k_c$  of the network of Figure 2 have an interpretation as statistical quantities, which follows from the derivation of Section 3.1. To maintain this interpretation, they have to satisfy certain relationships, namely (12), (16), (17), and (18). Some of these relationships do not necessarily need to be enforced. For example, since (18) is the only to involve  $\pi_c$ , there is a one to one relationship between the values of  $\log \kappa_c$  and  $\pi_c$ , independently of the value of  $|S_c|^{1/2}$ . Hence, it is equivalent to learn  $\pi_c$  under the constraint (18) or simply learn  $\log \kappa_c$ . The overall architecture of the MFAFVNet is shown in Since the latter leads to a simpler optimization, it should be

$$L(\mathcal{D}) = L_c(\mathcal{D}) + \lambda_1 \sum_c ||\Omega_c - P_c \Lambda_c||_F^2 + \lambda_2 \sum_c ||P_c - P_c^T||_F^2$$
(26)

where  $||A||_F$  is the Frobenius norm of A, the last two terms connected layer of appropriate dimensions and creates the are the regularizers that enforce the constraints and  $L_c(.)$  is input to the MFA-FV layer. This is implemented as shown a standard classification loss. In our implementation this is

$$L_{c}(\mathcal{D}) = \frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} \left[ \max(0, 1 - \delta(y_{i} = k) s_{k}(x_{i})) \right]^{p}$$
(2)

are finally included in the network before the final linear network,  $\delta(\cdot)$  the indicator function (1 when the argument holds and zero otherwise), and  $\lambda_1, \lambda_2$  two parameters used

nitialization MIT Indoor SUN

IIIIIaiizatioii	MIT IIIu001	SUN	
AlexNet			
Random	69.82	50.23	
Pre-Trained MFA	71.44	54.14	
VGG-16			
Random	77.3	56.2	
Pre-Trained MFA	80.3	62.51	

loss could in principle be used.

### 4. Experiments

In this section, we report on an extensive experimental evaluation of the MFAFVNet.

# 4.1. Experimental Setup

MIT Indoor scenes dataset [22] and the 397 class MIT SUN tialization of the resulting parameters is discussed below. dataset [28]. MIT Indoor includes 80 images of each category for training and 20 images for testing. SUN includes multiple train/test splits, with 50 images per class in the testing set. We present results for the average accuracy over tively. For both datasets, the complete network was trained

vious methods for scene classification. The VLAD of ization layer of the MFAFVNet was concatenated with the [11], the Sparse and H-Sparse coding of [17, 18], the Se-output of the penultimate layer of the Places network. mantic Fisher Vector (SFV) of [5], the full (FBN) and compact (compact BN) bilinear pooling networks of [7], 4.2. Parameter Initialization the MFA-FS of [6], the Deep FisherNet of [26] and the MetaClass method of [27]. While most of these methods [11, 5, 6, 26, 27] present results for both MIT Indoor and SUN, some [7, 18] only report in MIT Indoor. With the exception of the Deep FisherNet of [26], all these results are obtained by simply using features extracted from CNN layers, without any finetuning of the network. We simple restate their result. [26] did not address scene classification, only presenting results for object detection on PAS-CAL VOC 2012. We implemented the network as described in [26] and present its results on MIT Indoor and SUN.

mented with three different object recognition networks weaker on MIT Indoor and about 4% weaker on SUN. It is trained on ImageNet [4]: Alexnet [13], VGG-16, and VGG-19 [24]. The object class probability vectors produced by important to rely on an initialization with a strong statistical these networks, per  $l \times l$  patch, was converted to its natural interpretation.

Table 1. Effect of initialization on MFAFVNet classification accuracy. Table 2. Effect of regularization strength on MFAFVNet classification accuracy.

AlexNet					
$\lambda$	0.01	0.1	1	10	100
Accuracy	70.69	71.11	71.44	71.42	71.43
VGG-16					
		VGG-	16		
$\lambda$	0.01	VGG- 0.1	16	10	100

to control the strength of the regularization. In our implementation we use p = 2. The choice of the hinge loss tor p(x) of Figure 3. The PCA layer reduced this 1,000 is mostly for consistency with the Fisher vector literature, dimensional vector to the one with 500 dimensions, which which is mostly based on SVMs. Any other classification was used to compute the MFA-FV. Input images were resized, by making the smaller side 512-pixel long and maintaining the original aspect ratio. Three patch sizes,  $l \in$  $\{96, 128, 160\}$  were used, producing between 590 and 1000 patches per image.

The MFA contained 100 mixture components and a 10 dimension latent variable subspace. This produced a vector of  $500 \times 100 \times 10$  dimensions at the output of the MFA-FV layer. The parameters of the fully connected layer (FC9) Datasets: All experiments were based on the 67 class at the network output were initialized randomly. The ini-Layer FC9 was learned with a learning rate of 0.001 and all other layers with a learning rate of 0.00001. Momenon 10 epochs. As is costumary in the literature, some results are presented for the combination of the MFA-FV and the Places network [29], a network learned on the large Places Baselines: The MFAFVNet was compared to eight pre-scene dataset. In this case, the output of the L2 normal-

ir experiments have shown that a good initialization f the the PCA and MFA-FV layers can lead to substantial ains in classification accuracy. The PCA layer was initialzed by a PCA transformation learned from all patches at the output p(x) of Alexnet or VGG. The low dimensional vectors at the output of the PCA layer were then used to learn the MFA parameters with the EM algorithm of [8]. able 1 compares this initialization to one where all parameters are randomly initialized with a zero mean Gaussian istribution of standard deviation 0.01. The results in the table refer to a single patch size of 96 and  $\lambda_1 = \lambda_2 = 1$ , but we observed a similar behavior for other configurations. **Implementation Details:** The MFAFVNet was imple-The performance of the randomly initialized network is 2%

Table 3. Effect of patch scale on MFAFVNet classification accuracy. "3 scale" denotes the combination of three scales.

	AlexNe	et	VGG-16		VGG-19	
	MIT Indoor	SUN	MIT Indoor	SUN	MIT Indoor	SUN
$96 \times 96$	71.44	54.14	80.3	62.51	80.5	62.62
$128\times128$	71.4	54.03	78.44	61.47	79.29	61.48
$160 \times 160$	69.89	52.51	78.01	61.22	78.44	61.31
3 scales	75.01	57.15	81.12	64.51	82.66	64.59

### 4.3. Influence of Regularization

We next investigated the importance of regularization, b considering different values of  $\lambda_1$  and  $\lambda_2$  in (26). For simplicity, we considered only the case where  $\lambda_1 = \lambda_2 =$ which was also adopted in the remaining experiments. For small values of  $\lambda$  the network is free to learn a model that does not reflect the constraints of the MFA-FV, i.e. of weak statistical significance. For larger  $\lambda$  the parameters reflect the MFA constraints and the network has a stronger statistical significance. Table 2 presents results on MIT Indoor for different values of  $\lambda$  and a single scale patch with l = 96. There is an improvement of up to 1% when  $\lambda$  increases from 0.01 to 1 and performance stays approximately constant for larger  $\lambda$ . This shows that it is important to enforce the statistical significance of the parameters. In all subsequent experiments we have used the value of  $\lambda = 1$ .

### 4.4. Impact of Multiple Scales

Various recent works [5, 6, 17] have shown that is important. tant combine multiple patch scales, since objects of different sizes can be informative for scene classification. Table 3 summarizes the effect of patch sizes  $(l \in \{96, 128, 160\})$ on the classification accuracy of the MFAFVNet, for the hand, these are the only connections that allow the mixture three object recognition models. While scale  $96 \times 96$  outperforms the other two, results improve substantially when the three scales are combined. This confirms the previous observations of [5, 6, 17] for the benefits of multi-scale feature combination.

Table 4 compares the MFAFVNet to various previous these probabilities. scene classifiers based on an object recognition network trained on ImageNet [2, 6, 7, 17, 6, 26, 1]. The FV of [2] **4.6. Comparison to Scene Classifiers** uses both Alexnet and VGG-16 as CNN model and 10 patch scales. [17] extracts feature vectors at the output of the first fully connected layer, for a single patch size of l = 128, and uses them to derive a sparse coding based FV. [6] uses an MFA-based Fisher vector to pool the local features extracted from AlexNet, VGG-16, or VGG-19, but has no fine tuning. [26] simplifies the Fisher vector to allow end to end training of the network.

The MFAFVNet achieves state-of-the-art results on both datasets for both networks, even outperforming [2] which combines patches of ten different scales. The improved per-

formance over [11, 17, 2] is justified by the fact that these works rely on a FV derived from a Gaussian mixture of diagonal covariance and no fine tuning of the network. The improvements over the sparse coding techniques [17, 18] suggest that the MFA is a better model for the statistics f features learned by deep CNNs. Overall, the closest competitor to the MFAFVNet is [6], which also uses an MFA-FV but does not support end-to-end network finetuning. The MFAFVNet improves the results of this method by 1 - 2%, even though [6] concatenates Fisher vectors of different patch scales, to form a vector that is three times as long as that of the MFAFVNet.

Somewhat surprising is the improvement of 6% over the only other method that tried to train the network end-to-end [26]. We have found that this is due to their simplification of the FV, which includes removing the weights of the mixture components and the normalization terms in the denominators of the posterior probabilities of (21). A similar simplification of the MFAFVNet incurred losses of significant magnitude. Note that the simplification is computationally significant since, as can be seen in Figure 3, these are the only terms that require inputs from the other mixture components k in the softmax of the bottom branch. On the other components to interact with each other. From a statistical standpoint, in the absence of the normalization, the posterior probabilities are not even probabilities and the model looses coherence. Note that much of the computation of EM is aimed to get the posterior probabilities right, since **4.5. Comparison to Object-based Scene Classifiers** they determine the mixture assignments of the samples  $x_i$ . In fact, the E-step is mostly about getting good estimates of

9] trained a network with the same architecture as Alexnet or VGG for scene image classification directly from the Places scene dataset. This contains 2.4M scene images. Comparisons to this network test the effectiveness of representations such as the Fisher vector for transfer learning. Since the dimension of the feature vector extracted by the Places network is 4096, the dimension of the MFA-FS was reduced to this value in these experiments<sup>1</sup>. A comparison of the two approaches is presented in Table 5. Interest-

Table 4. Performance of sccene classifiers based on object recognition Table 6. Performance of combinations of object-based and scene-

Method	MIT Indoor	SUN	
AlexNet-based			
Sparse Coding [17]	68.2	-	
VLAD [11]	68.88	51.98	
FV [5]	72.86	54.4	
MFA-FS [6]	73.58	55.95	
FV+FC [2]	74.4	-	
MFAFVNet	75.01	57.15	
VGG-based			
Compact BN [7]	76.17	-	
Deep FisherNet [26]	76.48	57.91	
Full BN[7]	77.55	-	
Sparse Coding [18]	77.6	-	
H-Sparse [18]	79.5	-	
MFA-FS [6]	81.43	63.31	
FV+FC [2]	81.0	-	
MFAFVNet	82.66	64.59	

Table 5. Comparison to a scene classifier learned on Places.

Method	MIT Indoor	SUN	
AlexNet			
Places	68.24	54.3	
MFAFVNet	74.86	56.96	
VGG-16			
Places	79.47	61.32	
MFAFVNet	80.72	64.08	

ingly, the object-based network outperforms the Places network by a significant amount (up to 6%). This is likely due to the fact that scenes involve complicated combinations of objects, which may appear at different scales and poses. A In this work, we considered the transfer of a deep CNN network that is trained holistically, i.e. over the whole image, likely has difficulty in inferring these object cues. On classification. An embedded implementation of the MFAthe other hand, the combination of the object-based network FV was proposed. This enabled the design of a network and the pooling operation of the Fisher vector is basically architecture, the MFAFVNet, that can be trained in an end just learning the statistics of object appearances and some to end manner. The new architecture is based on a MFA-FV aspects of their configurations in the scene, e.g. relative layer that implements a statistically correct version of the properties such objects that appear at different sizes in a MFA-FV, through a combination of network computations class of scenes. In any case, these results show that objects and regularization. When compared to previous neural imare informative for scene classification. As far as we know, plementations of Fisher vectors, the MFAFVNet relies on this is also the only task on which transfer learning outper- a more powerful statistical model and a more accurate imforms the training of a deep network directly from a large plementation. When compared to previous non-embedded dataset of the target domain.

### 4.7. Combined networks

based models. This is, in fact, done in most previous works and regularization and the benefits of end to end training. and shown to improve performance. Following the standard practice in these experiments, we concatenate the feature vectors extracted by the two networks and classify and when combined with the holistic Places representation.

based scene classifiers.

Method	MIT Indoor	SUN		
AlexNet				
MetaClass+Places [27]	78.9	58.11		
FV+Places [5]	79.0	61.72		
MFA-FS+Places [6]	79.86	63.16		
MFAFVNet+Places	80.47	64.1		
VGG-16				
Deep FIsherNet+Places [26]	78.81	59.7		
MFA-FS+Places [6]	87.23	71.06		
MFAFVNet+Places	87.97	72.01		

the resulting vector with a linear SVM of hyperparameter  $C_{sum} = 2$ . These experiments did not use VGG-19, which has not been used in the Places network.

Table 6 shows the results obtained by combining the two networks. The combination of the Fisher vector with the holistic scene representations achieves a big improvement (up to 8%) over the performance of either of the representations independently, on both MIT Indoor and SUN. Since the networks capture complementary information the scene gist for Places and the object composition of the scene for the MFAFVNet - this suggests that these two classes of information are important, even complementary, for scene classification. Table 6 also shows that, even after combination with Places, the MFAFVNet achieves the best results among all Fisher vector representations. These results are, to the best of our knowledge, the state-of-art for scene image classification.

# 5. Conclusion

models, the MFAFVNet relies on a state of the art model, which is now embedded into a CNN. Experiments have shown the importance of maintaining a valid statistical in-It is also possible to combine the object- and scene-terpretation for the network, through proper initialization