

# Siamese Network of Deep Fisher-Vector Descriptors for Image Retrieval

Eng-Jon Ong, Sameed Husain and Miroslaw Bober  
University of Surrey  
Guildford, UK

e.ong, sameed.husain, m.bober@surrey.ac.uk

February 2, 2017

## Abstract

This paper addresses the problem of large scale image retrieval, with the aim of accurately ranking the similarity of a large number of images to a given query image. To achieve this, we propose a novel Siamese network. This network consists of two computational strands, each comprising of a CNN component followed by a Fisher vector component. The CNN component produces dense, deep convolutional descriptors that are then aggregated by the Fisher Vector method. Crucially, we propose to simultaneously learn both the CNN filter weights and Fisher Vector model parameters. This allows us to account for the evolving distribution of deep descriptors over the course of the learning process. We show that the proposed approach gives significant improvements over the state-of-the-art methods on the Oxford and Paris image retrieval datasets. Additionally, we provide a baseline performance measure for both these datasets with the inclusion of 1 million distractors.

## 1 Introduction

The rise of digital cameras and smart phones, the standardization of computers and multimedia formats, the ubiquity of data storage devices and the technological maturity of network infrastructure has exponentially increased the volumes of visual data available on-line and off-line. With this dramatic growth, the need for an effective and computationally efficient content search system

has become increasingly important. Given a large collection of images and videos, the aim is to retrieve individual images and video shots depicting instances of a user-specified object (query). There are a range of important applications for image retrieval including management of multimedia content, mobile commerce, surveillance, augmented automotive navigation etc. Performing robust and accurate visual search is challenging due to factors such as changing object viewpoints, scale, partial occlusions, varying backgrounds and imaging conditions. Additionally, today's systems must be highly scalable to accommodate the the huge volumes of multimedia data, which can comprise billions of images.

In order to overcome these challenges, a compact and discriminative image representation is required. Typically, this is achieved by the aggregation of multiple local descriptors from an image into a single high-dimensional global descriptor. The similarity of the visual content in two images is determined using a distance metric (e.g. Hamming or Euclidean distance) between their corresponding global descriptors. The retrieval is accomplished by calculating a ranking based on the distances between a set of images to a given query image.

This paper addresses the task of extracting a global descriptor by means of aggregating local deep descriptors. We achieve this using a novel combined CNN and Fisher Vector model that is learnt simultaneously. We also show our proposed model provides significant improvements in the retrieval accuracy when compared with related state-of-the-art approaches across different descriptor dimen-

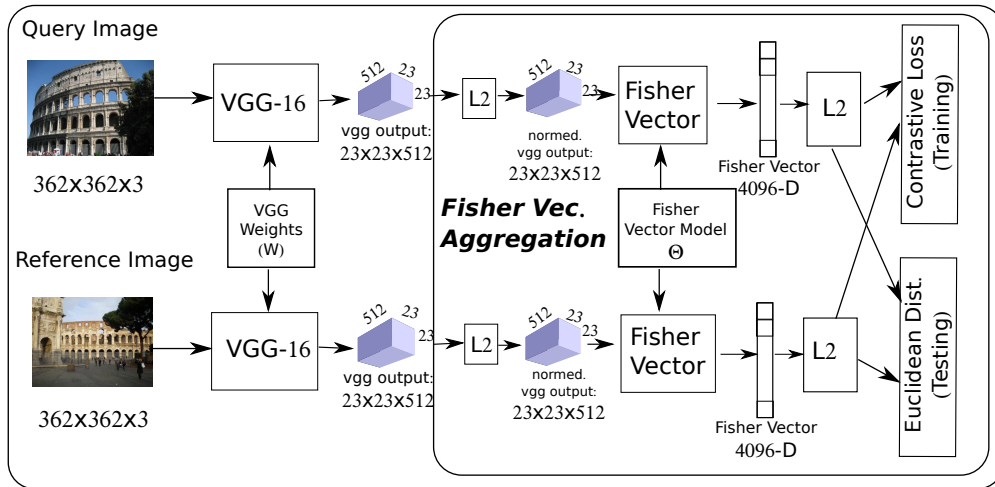


Figure 1: Overview of the training configuration of the proposed CNN-FV siamese network. During training, the last layer is the contrastive loss layer. During testing, the Euclidean distance between the final two fisher vectors is given.

sionalities and datasets.

## 1.1 Related Work

One popular method for generating global descriptors for image matching is the Fisher Vector (FV) method, which aggregates local image descriptors (e.g. SIFT [9]) based on the Fisher Kernel framework. A Gaussian Mixture Model (GMM) is used to model the distribution of local image descriptors, and the global descriptor for an image is obtained by computing and concatenating the gradients of the log-likelihoods with respect to the model parameters. One advantage of the FV approach is its encoding of higher order statistics, resulting in a more discriminative representation and hence better performance [11].

The FV model is learnt using unsupervised clustering, and therefore cannot make use of matching and non-matching labels that are available in image retrieval tasks. One way of overcoming this shortcoming was proposed by Perronnin et al. [12], where a fully connected neural network (NN) was trained by using the FV global descriptors as input. Here, the fisher-vector model was initially learnt in an unsupervised fashion on extracted SIFT features. The FV model then produces input feature vectors for the fully connected NN, which in turn is learnt in a

supervised manner using backpropagation.

However, both the SIFT features and FV model in the above method are unsupervised. An alternative is to replace the low-level SIFT-features with deep convolutional descriptors obtained from convolutional neural networks (CNNs) trained on large-scale datasets such as ImageNet. Recent research has shown that image descriptors computed using deep CNNs achieve state-of-the-art performance for image retrieval and classification tasks. Babenko et al. [2] aggregated deep convolutional descriptors to form global image representations: FV, Temb and SPoC. The SPoC signature is obtained by sum-pooling of the deep features. Razavian et al. [17] compute an image representation by the max pooling aggregation of the last convolutional layer. The retrieval performance was further improved when the RVD-W method was used for aggregation of CNN-based deep descriptors [6].

All of the above approaches use fixed pre-trained CNNs. However, these CNNs are trained for the purpose of image classification (e.g. 1000 classes of ImageNet) and may perform sub-optimally in the task of image retrieval. To tackle this, Radenovic et al. [16] and Gordo et al. [4] both proposed to use a Siamese CNN with max-pooling for aggregation. The CNN was fine-tuned on an image retrieval dataset. Two types of loss

function were considered for optimisation: 1) the contrastive loss function [16] and 2) the triplet loss function [4]. Both were able to achieve significant improvements from existing retrieval mAP scores. However, both these approaches use max-pooling as an aggregation method. The work proposed in this paper improves on this by employing a Fisher Vector model for aggregation instead of max-pooling. We also consider an alternative method of sum-pooling and compare different aggregation methods on standard benchmarks.

## 1.2 Contributions and Overview

The main contribution of this paper is a Siamese deep net that aggregates CNN-based local descriptors using the Fisher Vector model. Importantly, we propose to learn the parameters of the CNN and Fisher vectors simultaneously using stochastic gradient descent on the contrastive loss function. This allows us to adjust the Fisher vector model to account for changes in the distribution of the underlying CNN features as they are learnt on image retrieval datasets. We also show that our proposed method improves on the retrieval performance of the following state-of-the-art approaches: Siamese CNN with max-pooling[16] and Triplet loss with max-pooling [4]. We show that our approach achieves mAP scores that equal or improve on state of the art results for the Oxford (81.5%) and Paris datasets (82.5%). Importantly, this was achieved without any segmentation of images used in [4]. We also provide a new baseline of retrieval performance of our method when 1 million distractors are included into the test datasets.

The rest of the paper is structured as follows: Section 2 describes the proposed CNN-FV Siamese network used in this paper. The details for learning this network is given in Section 3. The experimental results are then described in Section 4 before concluding in Section 5.

## 2 Deep Fisher Vector Siamese Network

In this section, we describe the novel DNN that will learn a deep fisher vector representation by simultaneously learning the fisher-vector model components along with

the underlying convolutional filter weights in a Siamese network. The overview diagram of the proposed deep Siamese Fisher Vector network is shown in Fig. 1.

Traditionally, a Siamese network consists of two parallel branches in the network, where both branches share the same convolutional weights. One branch is fed a query image and the other branch a reference image which propagate through the network yielding 2 global descriptors respectively, which can be compared using Euclidean distance. Our proposed Siamese network is different in that each branch consists of two components: a CNN for producing deep image descriptors that are then aggregated via a Fisher Vector layer to produce the final global descriptor.

### 2.1 CNN-based Deep Descriptors

Suppose the input image is given as  $\mathbf{x} \in \mathcal{R}^{S \times S \times 3}$ . In order to extract the deep convolutional descriptors from the CNN component, the input image is first passed through number of convolutional layers. Here, we use convolutional layers with the same structure as the VGG-16 [18] network with the fully connected layers removed.

The CNN is effectively parameterised by the filter weights at each of its convolutional layers. We shall denote the collection of all the CNN filter weights as  $W$ . Formally, the CNN component can then be described by the function  $f : \mathcal{R}^{S \times S \times 3} \rightarrow \mathcal{R}^{O \times O \times F}$ , where  $F$  is the final number of convolutional filters, each producing a convolutional image of size  $O \times O$ . We then treat the final layer as producing a set of  $N_C = O \times O$  number of deep convolutional features that are of dimension  $F$ .

### 2.2 Fisher Vectors

In order to aggregate the  $N_C$  deep convolutional features, we employ the method of Fisher Vectors. Firstly, let  $\mathcal{X} = \{x_t \in \mathbb{R}^d, t = 1 \dots T\}$  be the set of  $N_C$   $F$ -dimensional deep convolutional features extracted from an image  $I$ . Let  $u_\Theta$  be an image-independent probability density function which models the generative process of  $\mathcal{X}$ , where  $\Theta$  represents the parameters of  $u_\Theta$ .

A Gaussian Mixture model (GMM) [13],  $u_\Theta$  is used to model the distribution of the convolutional features,

where:

$$u_{\Theta}(x) = \sum_{j=1}^C \omega_j u_j(x)$$

We represent the parameters of the  $C$ -component GMM by  $\Theta = (\omega_j, \mu_j, \Sigma_j : j = 1, \dots, C)$ , where  $\omega_j, \mu_j, \Sigma_j$  are respectively the weight, mean vector and covariance matrix of Gaussian  $j$ . The covariance matrix of each GMM component  $j$  is assumed to be diagonal and is denoted by  $\sigma_j^2$ . The GMM assigns each descriptor  $x_t$  to Gaussian  $j$  with the soft assignment weight ( $\tau_{tj}$ ) given by the posteriori probability:

$$\tau_{tj} = \frac{\exp(-\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu}_j)^T \Sigma_j^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_j))}{\sum_{i=1}^n \exp(-\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu}_i)^T \Sigma_i^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_i))} \quad (1)$$

The GMM can be interpreted as a probabilistic visual vocabulary, where each Gaussian forms a visual word or cluster. The  $d$ -dimensional derivative with respect to the mean  $\boldsymbol{\mu}_j$  of Gaussian  $j$  is denoted by  $\zeta_j$ :

$$\zeta_j = \frac{1}{T\sqrt{\omega_j}} \sum_{t=1}^T \tau_{tj} \Sigma_j^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_j) \quad (2)$$

We denote the elements of  $\zeta_j$  as  $\zeta_{jk}, k \in \{1, \dots, d\}$ . The final FV representation used,  $\zeta$ , of image  $I$  is obtained by concatenating the gradients  $\zeta_j$  for all Gaussians  $j = 1..n$  and normalising, giving:  $\hat{\zeta} = (\hat{\zeta}_{1,1}, \hat{\zeta}_{1,2}, \dots, \hat{\zeta}_{C,d})$ , with  $\hat{\zeta}_{jk} = \zeta_{jk}/|\zeta|$ , where  $|\zeta| = \sum_{j,k}^{C,d} \sqrt{\zeta_{jk}^2}$ . The dimensionality of  $\zeta$  is  $d \times C$ . Since the FV  $\zeta$  will be integrated into a Siamese-CNN, we shall henceforth refer to  $\zeta$  as ‘‘SIAM-FV’’ for **SIAM**ese-CNN-based **F**isher **V**ector.

### 2.3 Fisher Vector Partial Derivatives

In this section, the partial derivatives of the Fisher vector  $\zeta$  with respect to its underlying parameters ( $\Theta$ ) are given. These partial derivatives will be used for learning the proposed deep net. Firstly, we give the partial derivatives for the element ( $\zeta_{jk}$ ) of  $\zeta_j$  for some cluster  $j \in \{1, \dots, C\}$

and dimension,  $k \in \{1, \dots, d\}$ :

$$\frac{\partial \zeta_{jk}}{\partial \omega_j} = -\frac{1}{2T(\omega_j)^{3/2}} \sum_{t=1}^T \frac{\tau_{tj}(x_{tk} - \mu_{jk})}{\sigma_{jk}} \quad (3)$$

$$\frac{\partial \zeta_{jk}}{\partial \sigma_{jk}} = \frac{1}{T\sqrt{\omega_j}} \sum_{t=1}^T (x_{tk} - \mu_{jk}) \left[ \frac{\sigma_{jk} \frac{\partial \tau_{tk}}{\partial \sigma_{jk}} - \tau_{tk}}{\sigma_{jk}^2} \right] \quad (4)$$

$$\frac{\partial \zeta_{jk}}{\partial \mu_{jk}} = \frac{1}{T\sqrt{\omega_j}} \sum_{t=1}^T \frac{[(x_{tk} - \mu_{jk}) \frac{\partial \tau_{tj}}{\partial \mu_{jk}} - \tau_{tj}]}{\sigma_{jk}} \quad (5)$$

$$\frac{\partial \zeta_{jk}}{\partial x_{tk}} = \frac{1}{T\sigma_{jk}\sqrt{\omega_j}} \left[ (x_{tk} - \mu_{jk}) \frac{\partial \tau_{tj}}{\partial x_{tk}} + \tau_{tj} \right] \quad (6)$$

The partial derivatives of  $\tau_{tj}$  in the above equations are detailed in Appendix A. The equations Eq. 3 - 5 are used for calculating the gradients of the cluster prior, cluster mean and cluster standard deviation in the FV model. Eq. 6 is used to backpropagate errors to the filter weights in the CNN component. We find that the partial derivatives of the final normalised fisher vector elements  $\hat{\zeta}_{jk}$  all have the following form:

$$\frac{\partial \hat{\zeta}_{jk}}{\partial \phi} = \frac{1}{|\zeta|} \frac{\partial \zeta_{jk}}{\partial \phi} - \frac{\zeta_{jk}}{|\zeta|^3} \sum_{j=1}^C \zeta_{jk} \frac{\partial \zeta_{jk}}{\partial \phi} \quad (7)$$

In order to obtain the exact partial derivative of  $\hat{\zeta}_{jk}$  with respect to a particular parameter, we substitute  $\phi$  with this parameter, look up the corresponding equation in Eq. 3-6, and substitute it into Eq. 7 above.

## 3 Deep Learning of Fisher Vector Parameters

It is possible to learn the Fisher Vector GMM parameters using the EM algorithm on the deep convolutional features. However, this is an unsupervised method that does not make use of available labelling information. In order to tackle this shortcoming, we propose performing *supervised* learning of the GMM parameters. To this end, we treat the learning of the GMM parameters as part of learning process of a DNN.

For the purpose of learning, we are given a training dataset of  $T$  pairs of images, each image with resolution  $S \times S$ . Each pair of training images is associated

with a label, where 1 denotes matching images and 0 denotes non-matching images. We denote the training dataset as:  $\{(X_i, X'_i, Y_i)_{i=1}^T\}$ , where  $X_i, X'_i \in \mathcal{R}^{S \times S}$  and  $Y_i \in \{0, 1\}$ . The value of the labels of  $Y_i$  is 0 for matching examples and 1 for non-matching examples.

Next, we describe the contrastive loss [5] used for learning the proposed FV-CNN network. Firstly, Euclidean distance is used to measure the difference between two Fisher vectors:  $D(\zeta, \zeta') = \|\zeta - \zeta'\|$ .

Now, let the CNN weights be  $W$  and the set of all the Fisher Vector parameters  $\Omega$ . The loss function is defined as:

$$L(W, \Omega, Y_i, (\zeta_i, \zeta'_i)) = \frac{1}{2} Y_i (D(\zeta_i, \zeta'_i))^2 + \frac{1}{2} (1 - Y_i) (\max(0, \beta - D(\zeta_i, \zeta'_i)))^2 \quad (8)$$

where  $\beta$  is the heuristically determined margin parameter.

In order to optimise the GMM and cluster weight parameters of the Fisher vector,  $\Phi$ , the partial derivatives of  $L$  with respect to these respective parameters:  $\partial L / \partial \phi, \forall \phi \in \Theta$  are used. For conciseness, we will not write the arguments  $(\hat{\zeta}, \hat{\zeta}')$  when referring to the distance function  $D$ . So, using the chain rule on  $L$  gives:

$$\frac{\partial L}{\partial \phi} = \underbrace{[YD - (1 - Y) \max(0, \beta - D) \delta_{\beta - D > 0}]}_{\partial L / \partial D} \frac{\partial D}{\partial \phi} \quad (9)$$

The first backpropagated partial derivative  $\partial L / \partial D$  determines the amount of error present in the Fisher vectors of matching or non-matching pairs. The partial derivatives  $\partial D / \partial \phi$  allows us to adjust the FV model parameters and can similarly be derived using the chain rule, giving:

$$\begin{aligned} \frac{\partial D}{\partial \phi} &= \sum_{i=1}^D 2(\hat{\zeta}_i - \hat{\zeta}'_i) \left( \frac{\partial \hat{\zeta}_i}{\partial \phi} - \frac{\partial \hat{\zeta}'_i}{\partial \phi} \right) \\ &= \sum_{j=1}^C \sum_{k=1}^d 2(\hat{\zeta}_{jk} - \hat{\zeta}'_{jk}) \left( \frac{\partial \hat{\zeta}_{jk}}{\partial \phi} - \frac{\partial \hat{\zeta}'_{jk}}{\partial \phi} \right) \end{aligned} \quad (10)$$

where  $\zeta$  and  $\zeta'$  are the 2 input Fisher vectors to the distance function  $D$  and the partial derivatives of  $\hat{\zeta}_{jk}$  and  $\hat{\zeta}'_{jk}$  detailed in Section 2.3.

The parameters are then updated by adding the present value to their respective partial derivatives multiplied by the learning rate  $\alpha$ :  $\phi_{t+1} \leftarrow \phi_t + \alpha \partial L / \partial \phi$ .

## Updating CNN Weights

The updating of the CNN weights  $W$  is performed in a similar manner to the standard backpropagation, with the following difference: The gradients backpropagated from the contrastive loss and fisher layer is given by:  $\partial L / \partial D \times \partial D / \partial x_{tk}$  from Eq. 9 and 10, with the partial derivatives  $\partial \zeta_{jk} / \partial x_{tk}$  (Eq. 6) inserted in place of  $\partial \zeta_{jk} / \partial \phi$ . Since the CNN part is located below the Fisher vector layer, the above Fisher Vector gradients will then be propagated downwards to update the CNN weights  $W$ .

## 4 Experiments

For our experiments, the siamese network was learned on the Landmarks dataset used in [16]. Testing was performed on two independent datasets: Paris [15] and Oxford Buildings [14] with the mean average precision score reported. To test large scale retrieval, these datasets are combined with 1 million Flickr images [3], forming the Oxford1M and Paris1M dataset respectively. We followed the standard evaluation procedure and crop the query images of Oxford and Paris dataset, with the provided bounding box. The PCA transformation matrix is trained on the independent dataset to remove any bias.

### 4.1 Network Details

For the CNN component, the convolutional layers and respective filter weights of the VGG-16 network [18] was used. However, the max-pooling and ReLU layer at the final convolution layer was removed. 8 clusters was used for the Fisher vector GMM model, with their parameters initialised using the EM algorithm. This resulted in FV of dimensionality 4096. For retrieval purposes, we then perform PCA or LDA and whitening on the 4096-D Fisher vector, reducing dimensionalities to: 128D, 256D and 512D. In order to learn the PCA or LDA model, when the Oxford Buildings dataset is tested, the Paris dataset is used to build the PCA/LDA model, and vice versa. The contrastive loss margin parameter was set to  $\beta = 0.8$ . We set the learning rate equal to 0.001, weight decay 0.0005 and momentum 0.5. Training is performed to at most 30 epochs.

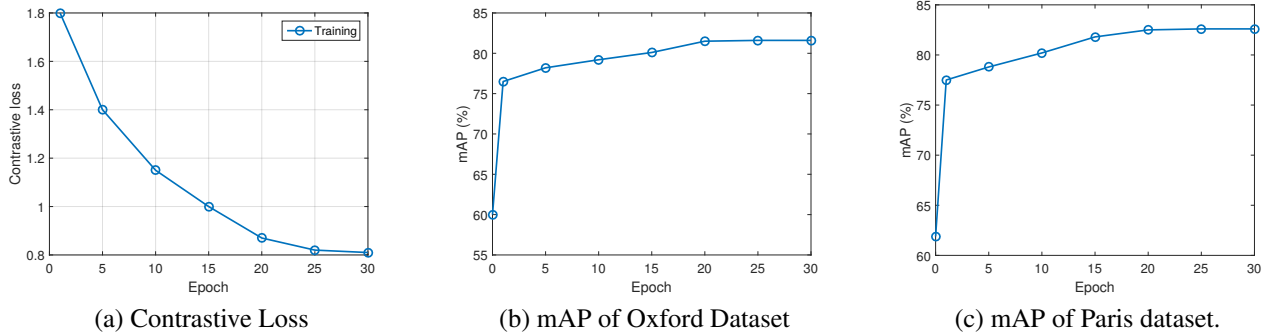


Figure 2: a) Contrastive loss value during the training of the proposed deep Siamese Fisher-Vector network across different epochs. b,c) show the mAP scores for the Oxford (b) and Paris (c) datasets across different training epochs.

## 4.2 Mining Non-Matching Examples

There exists significantly more non-matching pairs compared with matching pairs. Therefore, exhaustive use of non-matching pairs will create a large imbalance in the number of matching and non-matching pairs used for training. In order to tackle this, only a subset of non-matching examples are selected via mining, which allows the selection of only “hard” examples used in [16]. Here, for each matching pair of images used, 5 of the closest non-matching pairs to the query are used to form the non-matching pairs.

In this paper, 2000 matching pairs from the Landmarks dataset are randomly chosen. For each matching pair, 5 closest non-matching examples are then chosen, forming a 5-tuple, consisting of the following: query example; matching example; 5 non-matching examples. This forms a training set of  $2000 + 5 \times 2000 = 12000$  pairs. This set of 12K pairs will be re-mined after 2000 iterations. In total, each epoch in the training cycle consists of 6000 iterations.

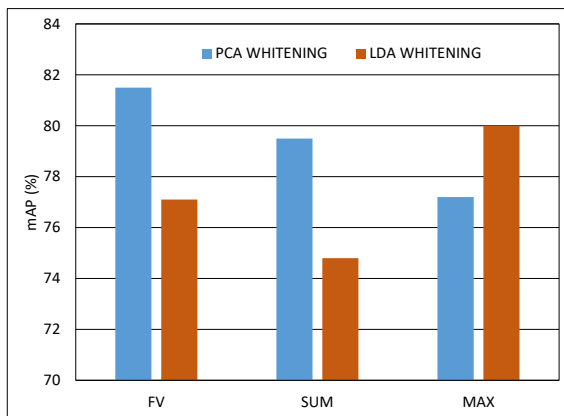
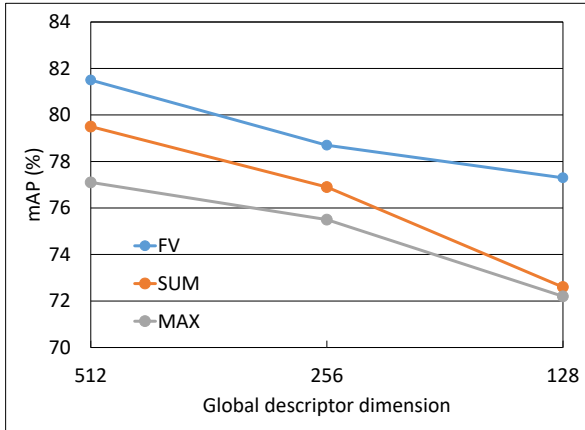


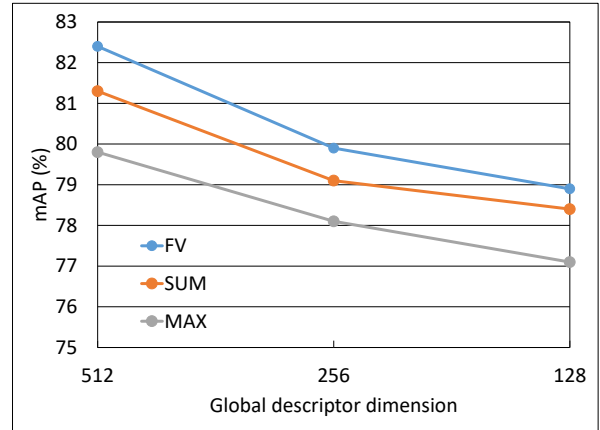
Figure 3: Shown are the mAP scores two dimensional-ity reduction methods: PCA and LDA on the proposed method (FV), sum-pooling and max-pooling [16] for the Oxford dataset.

## 4.3 Results

In this section, we evaluate the different components of our system in terms of: retrieval performance of the SIAM-FV descriptor across different epochs; projection methods (PCA vs LDA); dimensionality of the SIAM-FV descriptor; and compare the performance to the latest state-of-the-art algorithms.



(a) Oxford



(b) Paris

Figure 4: Shown are the mAP scores for the different global descriptor methods: proposed FV global descriptor (FV), sum-pooling (SUM) and max-pooling (MAX) [16], across different PCA-reduced dimensions and datasets: (a) Oxford; (b) Paris.

## Learning

The behaviour of the contrastive loss during learning is shown in Fig.2. It can be seen here that the initial 20 epochs give a large reduction in the loss value, and subsequent epochs producing only small further improvements in the loss function. In terms of the mAP results on the test datasets of Oxford and Paris, we find that the greatest improvement is obtained from the initial 5 epochs, with approximately 14-16% improvement in mAP scores. This can be seen in Fig. 2b) for the Oxford dataset and Fig. 2c) for the Paris dataset. Examples of the retrieved images based on the SIAM-FV descriptor for the Oxford and Paris datasets can be seen in Fig. 5 and 6 respectively.

## Projection Methods: PCA vs LDA

Fig. 3 shows the mAP results achieved by employing PCA and LDA for dimensionality reduction on the Oxford dataset. In [16], it was found that for max-pooling aggregation, LDA provided better performance at 80.0%, compared to PCA 76.1%. However, the converse was found for our SIAM-FV descriptor, which achieves 81.5% with PCA and 77.1% using LDA. This was also found to be the

case when sum-pooling was used, with 79.5% for PCA vs 74.8% using LDA. Thus for the remaining experiments, we have employed PCA as our choice for dimensionality reduction.

## Dimensionality of SIAM-FV

Figure. 4a,b, demonstrates the performance of SIAM-FV signature when reduced to different dimensionalities via PCA+Whitening. As expected, the best performance is obtained when the dimensionality is highest, at 512D for both Oxford and Paris datasets. Crucially, the proposed SIAM-FV has a mAP score that is approximately 2% higher than sum-pooling and 4% higher than max-pooling on the Oxford dataset across all dimensionalities 128D,256D and 512D. This gain in performance is similar for the Paris dataset, with the SIAM-FV method outperforming both sum-pooling and max-pooling across all dimensionalities.

## Comparison with State-of-the-Art

This section compares the performance of the proposed method to the state-of-the-art algorithms. Table 1 sum-

Table 1: Comparison with the state of the art using medium footprint signatures.

Method	Size	Oxf5k	Oxf105k	Paris6k
TEmb [7]	1024	56.0	50.2	-
NetVLAD [1]	4096	71.6	-	79.7
MAC [16]	512	58.3	49.2	72.6
R-MAC [19]	512	66.9	61.6	<b>83.0</b>
CroW [8]	512	68.2	63.2	79.7
MAC* [16]	512	80.0	75.1	82.9
SUM Pool	512	79.5	75.0	81.3
SIAM-FV	512	<b>81.5</b>	<b>76.6</b>	82.4

marises the results for medium footprint signatures (4k-512 dimensions). It can be seen that the proposed SIAM-FV representation outperforms most of the prior-art methods. On Paris dataset, the R-MAC representation provides marginally better performance. Note that R-MAC used region based pooling where deep features are max-pooled in several regions of an image using multi-scale grid.

Gordo et al. [4] achieved 83.1% on Oxford dataset. However they employed a region proposal network and extracted MAC signatures from 256 regions in an image, significantly increasing the extracting complexity of the representation.

We now focus on a comparison of compact representations which are practicable in large-scale retrieval, as presented in Table 2. The dimensionality of the SIAM-FV descriptor is reduced from 4096 to 128 via PCA. The results show that our method outperforms all presented methods. On the large dataset of Oxford1M SIAM-FV provides a gain of +2.4% compared to latest MAC\* signature.

## 5 Conclusions

In this paper, we have proposed a robust and discriminative image representation by aggregating deep descriptors using Fisher vectors. We have also proposed a novel learning method that allows us to simultaneously fine-tunes the deep descriptors and adapt the Fisher vector GMM model parameters accordingly. This effectively allows us to perform supervised learning of the Fisher vector model using matching and non-matching labels by op-

timising the contrastive loss. The result is a CNN-based Fisher vector (SIAM-FV) global descriptor. We have also found that PCA was a more suitable dimensionality reduction method compared with LDA when used with the SIAM-FV representation. We have shown that this model produces significant improvements in the retrieval mean average precision scores. On the large scale datasets, Oxford1M and Paris1M, SIAM-FV representation achieves a mAP of 62.5% and 63.2%, all yielding superior performance to the state-of-the-art.

## A Partial Derivatives of $\tau_{tj}$

We find that the partial derivatives of  $\tau_{tj}$  with respect to  $\sigma_{jk}$  and  $\mu_{jk}$  both have the same form. So, let  $\phi$  be either  $\sigma_{jk}, \mu_{jk}$  or  $x_{tk}$ . Also, let the numerator of  $\tau_{tj}$  be denoted as  $\tau_{tj}^{(j)}$  and its denominator  $\tau_{tj}^{(\Sigma)}$ , so  $\tau_{tj} = \tau_{tj}^{(j)} / \tau_{tj}^{(\Sigma)}$ , then:

$$\frac{\partial \tau_{tj}}{\partial \phi} = \frac{[\tau_{tj}^{(\Sigma)} - \tau_{tj}^{(j)}] \frac{\partial \tau_{tj}^{(j)}}{\partial \phi}}{(\tau_{tj}^{(\Sigma)})^2}, \quad \phi \in \{\sigma_{jk}, \mu_{jk}, x_{tk}\}$$

where,

$$\begin{aligned} \frac{\partial \tau_{tj}^{(j)}}{\partial \mu_{jk}} &= \tau_{tj}^{(j)} \left[ \frac{(x_{tk} - \mu_{jk})^2}{\sigma_{jk}^3} \right] \\ \frac{\partial \tau_{tj}^{(j)}}{\partial \sigma_{jk}} &= \tau_{tj}^{(j)} \left[ \frac{(x_{tk} - \mu_{jk})}{\sigma_{jk}^2} \right] \\ \frac{\partial \tau_{tj}^{(j)}}{\partial x_{tk}} &= \tau_{tj}^{(j)} \left[ \frac{\mu_{jk} - (x_{tk})}{\sigma_{jk}^2} \right] \end{aligned}$$

## References

- [1] R. Arandjelović, P. Gronat, A. Torii, T. Pajdla, and J. Sivic. NetVLAD: CNN architecture for weakly supervised place recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [2] A. Babenko and V. S. Lempitsky. Aggregating deep convolutional features for image retrieval. *CoRR*, 2015.
- [3] M. Bober, S. Husain, S. Paschalakis, and K. Wnukowicz. Improvements to TM6.0 with



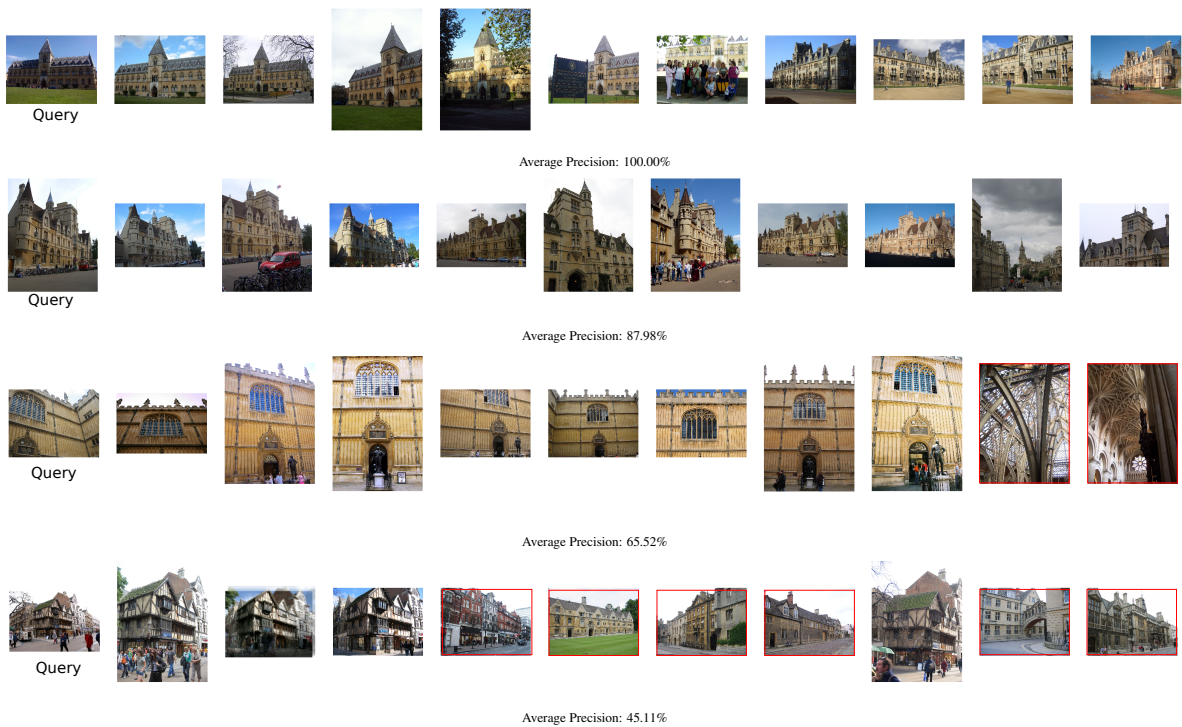


Figure 5: Examples of the top 10 retrieved images on the Oxford dataset using the proposed method for different average precisions. Retrieved non-matching images are highlighted in red.

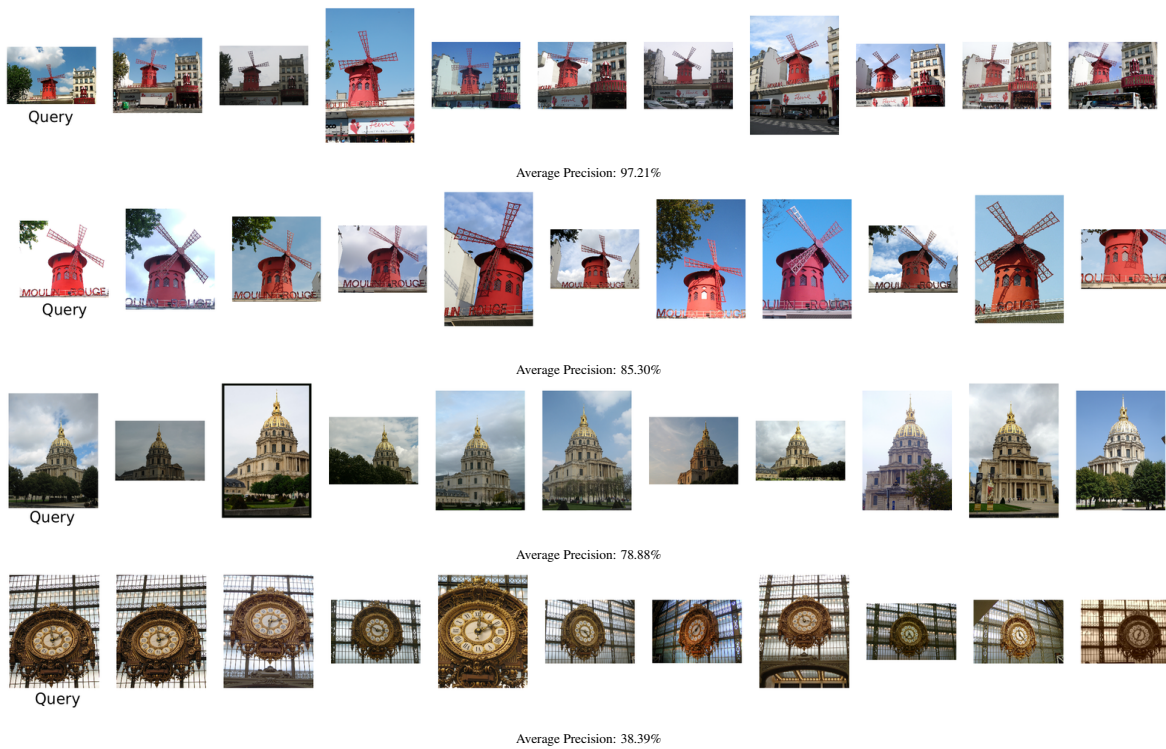


Figure 6: Examples of the top 10 retrieved images on the Paris dataset using the proposed method for different average precisions.

Table 2: Comparison with the state of the art using small signatures.

Method	Size	Oxf5k	Oxf105k	Oxf1M	Paris6k	Paris1M
Max-pooling [17]	256	53.3	-	-	67.0	-
SPoC [2]	256	53.1	50.1	-	-	-
MAC [16]	256	56.9	47.8	-	72.4	-
NetVLAD [1]	256	63.5	-	-	73.5	-
CroW [8]	256	65.4	59.3	-	77.9	-
Ng et al [10]	128	59.3	-	-	59.0	-
MAC* [16]	128	76.8	70.8	60.1	78.8	62.5
SUM Pool	128	72.6	67.7	57.9	78.4	62.4
SIAM-FV	128	<b>77.3</b>	<b>71.8</b>	<b>62.5</b>	<b>78.9</b>	<b>63.2</b>

- a robust visual descriptor proposal from University of Surrey and Visual Atoms. In *MPEG Standardisation contribution : ISO/IEC JTC1/SC29/WG11 CODING OF MOVING PICTURES AND AUDIO, M30311*, jul 2013.
- [4] A. Gordo, J. Almazán, J. Revaud, and D. Larlus. *Deep Image Retrieval: Learning Global Representations for Image Search*, pages 241–257. IEEE Computer Society, 2016.
- [5] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *Proc. of CVPR 2006, CVPR '06*, pages 1735–1742, Washington, DC, USA, 2006. IEEE Computer Society.
- [6] S. S. Husain and M. Bober. Improving large-scale image retrieval through robust aggregation of local descriptors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016.
- [7] H. Jégou and A. Zisserman. Triangulation embedding and democratic aggregation for image search. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- [8] Y. Kalantidis, C. Mellina, and S. Osindero. Cross-dimensional weighting for aggregated deep convolutional features. *CoRR*, 2015.
- [9] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, pages 91–110, 2004.
- [10] J. Y. H. Ng, F. Yang, and L. S. Davis. Exploiting local features from deep networks for image retrieval. In *2015 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 53–61, 2015.
- [11] F. Perronnin and C. R. Dance. Fisher kernels on visual vocabularies for image categorization. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2007.
- [12] F. Perronnin and D. Larlus. Fisher vectors meet neural networks: A hybrid classification architecture. In *Proc. of CVPR*, pages 3743–3752. IEEE Computer Society, 2015.
- [13] F. Perronnin, Y. Liu, J. Sanchez, and H. Poirier. Large-scale image retrieval with compressed fisher vectors. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3384–3391, 2010.
- [14] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman. Object retrieval with large vocabularies and fast spatial matching. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2007.
- [15] J. Philbin, M. Isard, J. Sivic, and A. Zisserman. Lost in quantization: Improving particular object retrieval in large scale image databases. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [16] F. Radenović, G. Tolias, and O. Chum. *CNN Image Retrieval Learns from BoW: Unsupervised Fine-Tuning with Hard Examples*, pages 3–20. IEEE Computer Society, 2016.
- [17] A. S. Razavian, J. Sullivan, A. Maki, and S. Carlsson. Visual instance retrieval with deep convolutional networks. *CoRR*, 2014.

- [18] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [19] G. Tolas, R. Sivic, and H. Jégou. Particular object retrieval with integral max-pooling of CNN activations. *CoRR*, 2015.