# Object Instance Identification with Fully Convolutional Networks

Maxime Portaz · Matthias Kohl · Jean-Pierre Chevallet · Georges Quénot · Philippe Mulhem

Received: date / Accepted: date

Abstract This paper presents a novel approach for instance search and object detection, applied to museum visits. This approach relies on fully convolutional networks (FCN) to obtain region proposals and object representation. Our proposal consists in four steps: a classical convolutional network is first fined-tuned as classifier over the dataset, next we build from this network a second one, fully convolutional, trained as classifier, that focuses on all regions of the corpus images, this network is used in a third step to define image global descriptors in a siamese architecture using triplets of images, and eventually these descriptors are then used for retrieval using classical scalar product between vectors.

Our framework has the following features: i) it is well suited for small datasets with low objects variability as we use transfer learning, ii) it does not require any additional component in the network as we rely on classical (i.e. not fully convolutional) and fully convolutional networks, and iii) it does not need region annotations in the dataset as it deals with regions in a unsupervised way. Through multiple experiments on two image datasets taken from museum visits, we detail the effect of each parameter, and we show that the descriptors obtained using our proposed network outperform those from previous stateof-the-art approaches.

**Keywords** Fully Convolutional Network  $\cdot$  Triplet Loss  $\cdot$  Siamese Network  $\cdot$  Instance Search  $\cdot$  Image Retrieval

# 1 Introduction

The work presented here is dedicated to enhance a museum audio-tour guide with a camera, in order to help user orientation, enable automatic guidance

Univ. Grenoble Alpes, CNRS, Grenoble-INP, LIG, F-38000 Grenoble France E-mail: Firstname.Lastname@imag.fr

and facilitate museum artifact explanations: when the visitor is close enough to an object, a multimedia explanation is automatically presented. The camera is used for the system to localize the user without the need of any extra hardware to be installed in the museum. Obviously, the entire museum must be photographed (or video recorded), and each object image then has to be localized in the digital museum map.

Instance search is a visual task that aims, given an image, to identify the particular objects shown. Instance search must not be confused with Image Classification that focuses on identifying object category, with robustness to intra-class variability. In addition, an instance (representing one given object) is generally described only by a few shots. The main difference between Image Classification and image retrieval, is the amount of data and their variability. In classification, we rely on a large amount of data with high variability of examples, and we can then train a Deep Convolutional Network with millions of parameters. In image retrieval or identification, as we want to identify a particular instance, the variability of examples is less important, and not sufficient to train a network like ResNet. In order to use a CNN, we only fine-tune a CNN pre-trained on a bigger collection for image classification.

In our work, the solution chosen to identify instances, using an image retrieval system is in two steps: in a first step we retrieve all images similar to a query image, and then we decide, from the image retrieval system result list, the instance identified.

We propose an image retrieval system that learns image representations with Deep Learning Convolutional Neural Networks (CNN). The Neural Network model proposed is learned with a siamese network with three streams and a triplet loss [25]. The aim is to produce an image representation that allows image comparison based on their contents. Because of the relatively small number of images available in Instance Search dataset, we need an external source of data to train a convolutional network, such as ImageNet [23]. The network we use is a Fully Convolutional Network (FCN) [15] that allows any input size, to avoid image deformation or scaling. One main reason that led us to use such networks is that the FCN can be used to produce region proposals without any additional component, in the network or in the dataset. Such property is needed in our case, as we want to be as prices as possible, even when the camera is not correctly positioned on the chest of the user. For the training phase, we use the triplet loss between the three streams of the siamese architecture, and a cross entropy loss for classification of the region with the highest activation. The aim is to create a representation of the image that captures the position of the object and the difference between images, whether similar or not. At test time, the trained FCN is passed over the whole image, but only the location with the top k maximal activations will form the image description. This representation is compared with the reference images in the dataset using a dot product, to obtain a ranked list of images in the deceasing order of similarity. Then, the closest reference image representing the instance is selected as the identified object.

We evaluate our approach on two egocentric datasets from museum visits [21]. We show that our approach achieves better results on these datasets than the previous state of the art by Gordo et al. [11].

In the following, the section 2 presents the related work on instance identification. Then, in section 3 we describe why the use of pre-trained CNNs and fine-tuning is important for our problem. The section 4 introduces the region of interest detection and object localization. The section 5 describes the proposed network and how it was trained on the datasets we used. Datasets used for evaluation are presented in section 6. Experimental results and evaluation are shown in section 7, in which we detail the impact of the parameters of our proposal on the quality of the results. We conclude in section 8 by giving some future directions of this work.

## 2 Related Work

Before the ground-breaking results of deep learning methods for object detection and image retrieval, shallow patch descriptors have been used in several domains. The SIFT [16], Scale-Invariant Feature Transform, descriptor was the most used one, among the large variety of traditional patch descriptors. It has been successfully employed for tasks like image search with content-based retrieval [13] or classification [18]. For image retrieval, methods inspired by text retrieval methods, such as bag of words [6], use bag-of-features (BoF) image representations [28] that group similar features together in clusters and stores the number of occurrences of each cluster in one image.

In order to compare images for image retrieval, image patch comparisons have shown better results than SIFT [8,27,32]. Image patches can be constructed with deep patch descriptors [8] as patch label, each patch is a label, by learning patch differences with a siamese network [27,32], or with a Convolutional Kernel Network [17].

Starting with the results of AlexNet for image classification in the 2012 ImageNet challenge [14,23], image classification tasks have been dominated by CNNs. A CNN trained on a large enough labeled dataset like ImageNet can be used as a feature extractor with its intermediate layers, to construct an image representation for image retrieval [4,26]. To overcome the lack of geometry invariance of this approach [9], cross-matching [26], sum-polling [3] or fine-tuning [4] with an external dataset can be used. Fine-tuning focuses on the higher layers of a CNN and can increase generalization even in the finetuned model[31]. These approaches are well suited for classification purpose, i.e., when we have numerous samples per class, but cannot be applied directly to instance search or image retrieval.

Another important aspect of image retrieval is to learn to rank [1,10]. While Arandjelovic et al. [1] have shown the importance of learning to rank, Gordo et al. [10] used a siamese network [7] along with a triplet loss, previously used for face recognition [25], to construct an effective image representation



**Fig. 1** Siamese Architecture use to train a network (NN) with a triplet of images, with an anchor image (a), a negative image (n) and a positive image (p)

by learning with a similarity metric.

As shown in figure 1, a three-way siamese architecture use a triplet of images  $\langle I_a, I_n, I_p \rangle$  with  $I_a$  being the anchor image (or query),  $I_n$  the negative example and  $I_p$  the positive one. The triplet loss is define to maximize the distance between the representation of the anchor  $(x^a)$  and the negative example one  $(x^n)$ , and minimize the one between anchor and the positive example  $(x^p)$ .

Previous approaches in image retrieval [11,24,29] usually deal with regions of interest in one way or another. The idea is that in most cases, only certain parts of each image can be useful for comparison with other images. In addition to this, cropping images at their regions of interest can help with differences in scale of the images to compare: if a painting is visible only in a small part of an image, cropping the image at that part and then re-scaling the part should set the painting at a normalized scale. However, in instance search with museum datasets, it is not obvious where the regions of interest should be: most images represent an entire painting or parts of it and only some may contain the painting as part of the image with a wall in the background. This means for most images, the ground-truth region of interest is simply the entire image, and some may have a ground-truth region of interest which is almost the entire image, excluding only a small part of the background.

As noted by previous authors, when using the triplet loss, it is crucial to choose the best triplets during training in order to obtain convergence. In particular, many triplets are irrelevant and do not produce any loss since they are *too easy* for the network.

Hence, the first idea is to choose the hardest triplets. However, this can lead to a collapsing model with a bad local minima early on in training, as explained by Schroff et al [25]. Thus, they choose *semi-hard* triplets instead. Semi-hard triplets are obtained as follows: use all possible positive couples of images (couples of images from the same instance). For each positive couple, choose the hardest negative that is easier than the positive couple. Hard and easy are defined by the dot product between the descriptors of the images: a high value of the dot product for images of the same instance represents an easy positive couple, a high value of the dot product for images of different instances represents a hard negative couple. The value of all dot products are determined before each pass over the whole training data during training, for all couples of images.

A different triplet selection mechanism was proposed by Gordo et al. [11]. First, calculate the values of dot products for all couples of images before each pass over the training data. Second, for each image, choose the n easiest positive images and the m hardest negatives. Then, calculate the loss for all possible combinations and use the o triplets with the highest loss. This method probably eliminates some noise when choosing the easiest positive couples, for images that are labeled as being the same instance but are not visually similar.

## 3 Fine-tuning and object Localization

The modularity of a CNN means that we can easily transfer the weights from a pre-trained model, and only re-train the highest abstraction layers. Specifically, we re-train all fully connected layers and the highest level convolutional layers in the model (depending on the architecture), since our datasets contain many visually different images as compared to the ImageNet dataset used for pre-training the models.

A network fine-tuned on classification on such a dataset should be able to easily identify the region containing the painting, since the background is contained in almost all classes, which means it is a particularly bad indicator of the class. Thus, if the network is applied in a strided manner across an image, it should produce low maximal activations in parts containing big sections of background wall.

Figure 2 shows images, along with the heat map representing the maximal activation of a fine-tuned ResNet-152 at each coordinate, when the network is applied in a strided manner across the input image. From this image, we can see that the highest maximal activations of the network usually occur at the location of the object. This is true even if the object is not correctly classified by some of the highest activations as can be seen in the second image.

In the third image, it seems like many high maximal activations occur specifically in the background area. However, the corresponding label-map shows that these areas correspond to the labels 38E and 43D. Both of these labels are pieces of art which consist mostly of the background wall. In this sense, it is not entirely wrong to consider 'wall-only' patches of the image as instances of these pieces of art. This simply means that the image consists of two separate regions of interest: one region with the painting (label 30P) and one region with the wall (labels 38E/43D).

From these observations, we can confirm the assumption that the maximal activations of a fine-tuned network are a good indicator of the location of an object, or a combination of different objects. Using this assumption, there is



Fig. 2 Sample images (scaled to a smaller side of 448 pixels) along with the heat-map of maximal activation values at each location when a fine-tuned ResNet-152 is applied to the image in a strided manner, as well as the labels of all maximal activations that are greater than the mean maximal activation

no need for a procedure to annotate regions of interest, as employed by most state-of-the-art image retrieval approaches [10,29,22].

On the other hand, using datasets developed for image retrieval, such as Paris6k or Oxford5k [20,19], this assumption cannot be applied, since the dataset is not clean enough for a network fine-tuned on classification to be a good indicator of location of the query objects.

## 4 Fine-tuning on classification using FCN

As shown before, a fine-tuned CNN is already a good indicator of the location of an object in our datasets. Additionally, it seems like scale is a particularly important factor.

Thus, the idea is to start by fine-tuning a network with images at different scales. This can be achieved by using a fully convolutional network (FCN) [15].

In an FCN, the final fully connected layers of a network are replaced by convolutional layers having a kernel which fits the entire domain of the output of the previous layer. This type of convolution is equivalent to a fully connected layer, but allows inputs (and outputs) of any size. The effect is that the network can be applied in one pass to an arbitrarily sized image. The output then represents the activations of the network as if it was applied in a strided manner across the image.

Once an FCN is applied to the image, the loss is calculated by averaging the cross-entropy (CE) loss (fig. 3). Given a scale s of an image, the loss  $L_s$  is computed by averaging the cross entropy(eq. 1) of every regions. The equation 1 show the equation for the cross-entropy loss, given an input x, that contains the score for each class and c the correct class.

The final loss is then obtained by passing images at different scales through the FCN and averaging across all cross-entropy losses of all outputs and scales.

$$\mathcal{CE}(x,y) = -y\log(softmax(x)) \tag{1}$$

We choose to give each scale of the image the same weight in the loss. This is because the images are passed to the network at their true aspect ratio, which means the loss for different images may have different values for the heights and widths of the feature maps  $H_s$  and  $W_s$ .

## 5 Constructing the image descriptor

#### 5.1 Training with siamese architecture

When training, the network is applied to a triplet of images (fig. 4) with a siamese configuration. The overall loss  $\mathcal{L}$  used for this training is the triplet loss defined in [25]. The triplet loss, defined initially in [25] using squared distances, can be expressed using dot products when considering normalized vectors. This leads to simpler gradient computation.

In experiments, we found that method developed by Gordo et al. [11] for triplets selection does not perform well for datasets with few images per instance, since we either have to choose n as very low or we end up choosing all positive couples after all for most instances, just like in the *semi-hard* selection. We choose the semi-hard triplet selection for the first two passes over the dataset, after which we only choose the hardest negatives for all positive couples.



Fig. 3 Loss computation when training the network over different regions and scales of the image. LSM is logsoftmax and NLL is Negative Log-Likelihood, there are used for cross entropy loss(eq. 1)



Fig. 4 Proposed architecture for instance search based on an FCN [15] for region proposals, at training time

Equation 2 shows the loss as used in our experiments to train the proposed model, for N images. In this equation,  $(h_l, w_l)$  represent the spatial coordinates of the *l*-th region of highest maximal activation in the feature map produced by the FCN.  $x_i^a x_i^n$  corresponds to the dot product between the *i* anchor descriptor and the *i* negative example descriptor and  $x_i^a x_i^p$  to the dot product between the anchor descriptor and the positive example descriptor. The scalar *m* represents the margin between a positive and a negative pair of images.

We regularize the triplet loss by a cross-entropy loss to make sure that the k locations with highest maximal activations are correctly classified. This loss is averaged over the k locations.

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^{N} \left( \max(0, x_i^a x_i^n - x_i^a x_i^p + m) + \alpha \frac{1}{k} \sum_{l=1}^{k} \mathcal{CE}_i^{h_l, w_l}(x_i^a, y_i^a) \right)$$
(2)



Fig. 5 Proposed architecture for instance search based on an FCN [15] for region proposals, at deploy time

In our experiments, we choose the number of regions with highest maximal activation to be k = 6 and the regularization hyper-parameter  $\alpha = 1.0$ . The margin of the triplet loss is m = 0.1.

This approach allows the network to decide which region of interest is best suited for classification and ultimately which regions are best suited for comparison with other images. Another advantage is that this approach does not require any annotation of the images with regions of interest, which can be a long, manual or automatic process, as evident from the cleaning process used by Gordo et al [11].

#### 5.2 Descriptor Extraction Network

Figures 5 illustrate the proposed architecture for image descriptor extraction. To obtain a descriptor, we first apply the convolutional layers of a previous architecture. We then obtain all classification outputs at all locations using the FCN. We only consider the maximal activation at all locations. The locations with the top k maximal activations will form the descriptor.

For each of these locations, the convolutional features are reduced by a  $\|\cdot\|_2$ -normalization, then a shifting and fully connected layer. Finally, all descriptors from the k locations are sum-aggregated and  $\|\cdot\|_2$ -normalized again.

An important property of the descriptor is that it heavily relies on the classification capabilities of the network. This means the descriptor is mostly meaningless for a different dataset and needs to be learned for each dataset. This can be an advantage, since the descriptor can be better suited to a particular dataset and the learning process does not take long. On the other hand, it means that the descriptor cannot be applied in a typical image retrieval task.

#### 5.3 Instance Feature Augmentation

An approach called Database-side feature augmentation [30,2], proposes to combine descriptors of the reference images in order to form better databaseside descriptors. Every reference descriptor is simply replaced by a combination of itself and the k nearest neighbors. This combination is computed as a weighted sum, weighted by the rank of the neighbors with respect to k (the closest neighbor has the highest weight and the k-th neighbor the lowest). In our work, we use a technique called Instance Feature Augmentation. We use the fact that we know the corresponding label for each image in our dataset. For each label, we compute the representation of an instance by averaging the features of every images corresponding to this label. This representation is added to the dataset as a new instance. We show that this approach does not improve mean precision@1, but gives a better Mean Average Precision. This suggests that the internal representation of the instance is improved.

# 6 Dataset

The proposed approaches as well as several baselines are evaluated on two datasets of still images, namely CLICIDE and GaRoFou. These datasets are described in detail by Portaz et al. [21]. They represent artwork photos, taken by classical or head-mounted cameras. Table 1 details their characteristics.

Corpus	#instances	#images	#queries	images/instance
CLICIDE	473	3425	177	7.24
GaRoFou	311	1252	184	4.03

Table 1 Characteristics of the two corpus considered in the experiments.

Both datasets are typical of instance search datasets in museums or touristic sites. The objects represented by their images are paintings for one, first column) and glass cabinets containing sculptures and artifacts for the other dataset (fig. 6). Both datasets contain a small number of images per instance and a small number of images in total, with respectively 4 and 6 images per instance in average.

## 7 Evaluation

# 7.1 Fine-tuning

In our experiments, we focus on two well studied networks: AlexNet [14] and ResNet [12].

#### 7.1.1 Layers selection

The experiments conducted in [21] on those datasets, show that the best results are obtained when fine-tuning the last convolutional layer and above. We conducted experiments by first retraining the last layer, and after few epochs and stabilization, add the previous last layer<sup>\*</sup>, and so on. This led to the following choices:

 For the AlexNet architecture, we choose to re-train all layers above and including the last convolutional layer.











Fig. 6 CLICIDE and GaRoFou Dataset example. The first column are images from CLI-CIDE dataset, representing painting. The right column are images from GaRoFou, with pictures of 3D objects, like sculptures.

 For a ResNet architecture, we re-train all layers above and including the third to last block of convolutional layers. This contains the nine highest convolutional layers in total.

This can be explained by the high specialization on the dataset of the last layers of the network. On the other side, a large amount of data is required to retrain deeper layers.

### 7.1.2 Data augmentation

Image retrieval methods focus on problems with few examples and little variability in instance images. This leads to too few data to train a typical CNN model designed for classification, even with fine-tuning. One way to overcome this is to augment the data, by randomly applying affine transformations, color perturbations and other random transformations.

The lack of geometry invariance and scaling invariance of the model can be reduced by randomly rotating and flipping the images and using different scaling, thus we perform this type of data augmentation throughout our experiments.

For data augmentation in order to fine-tune a CNN, we use the following values in our experiments:

- 1. Rotation: any angle is chosen with the same probability.
- 2. Scaling: the scaling factor is chosen independently for each dimension in the range [0.75, 1.25].
- 3. Flipping: with probability 0.5, images are horizontally flipped.

#### 7.2 Parameter for the Fully Convolutional Network

The stride of a full network depends on the architecture and is 32 pixels for the architectures used here: AlexNet and ResNet.

For the processing of the Fully Convolutional networks (step 2 of our proposal, described in part 5), all images are scaled to have the same number of pixels in the smaller side in order to normalize the sizes of the features present in the images. Note that for large aspect ratios and large scales of the smaller side, the memory consumption of training can be high for single images having a very large aspect ratio. To limit this spike in memory consumption, the aspect ratios are limited by introducing uniform random noise on the smaller side of images with high aspect ratios. In our experiments, we use a maximal aspect ratio of 2.0 and images at two scales of 448 and 224 pixels for the smaller side. We found that the AlexNet architecture did not have good convergence behavior, thus we used scales of 384 and 224 instead.

#### 7.3 Results

Table 3 gives an overview of the results obtained. First, the baselines established by SIFT descriptor, CNN network features extraction are shown. Addi-

	Mean Precision@1 (in %)
	CLICIDE
Proposed AlexNet FT only	79.39
$Proposed \ AlexNet \ FT + FT$ -region	81.21
Proposed ResNet-152 FT only	92.73
Proposed ResNet-152 $FT + FT$ -region	94.55

Table 2 Evaluation of the influence of region fine-tuning on the final model.

	Mean Precision@1 (in %)		Mean Ave. Precision (in %)	
	CLICIDE	GaRoFou	CLICIDE	GaRoFou
SIFT [21]	70.08	78.82	N/A	N/A
ResNet-50 [10]	90.30	95.65	65.49	88.43
ResNet-50, multi-res [10]	92.73	95.65	N/A	89.32
AlexNet IN	72.73	85.87	32.71	66.11
AlexNet FT	78.18	90.76	38.51	72.92
AlexNet SS	75.76	90.20	36.20	77.73
Proposed AlexNet	81.21	83.15	45.53	71.71
Proposed AlexNet (IFA)	80.61	82.61	71.02	81.66
ResNet-152 IN	72.12	85.33	40.99	70.15
ResNet-152 FT	79.39	94.57	75.11	93.44
ResNet-152 SS	85.45	95.11	83.00	91.90
Proposed ResNet-152	94.55	96.20	82.94	91.83
Proposed ResNet-152 (IFA)	93.94	95.11	94.23	93.86

 Table 3
 Mean precision@1 and mean average precision evaluation results for the CLICIDE and GaRoFou datasets.

tionally, we show the relevant results obtained by fine-tuning a classification network, abbreviated by FT in the table. We then show the results obtained by a simplified Siamese architecture, abbreviated SS. Finally, we show the results obtained by the proposed network. In addition to the mean precision@1, we show the mean average precision obtained by the different approaches.

From the baselines presented, we can make two observations. First, even a simple global descriptor obtained from the convolutional features of a CNN pre-trained on ImageNet performs better than matching local SIFT descriptors on our datasets. Second, the ResNet-50 proposed by Gordo et al. [10] outperforms the descriptors from pre-trained networks by far, even though it has never seen the images from our datasets during training, either.

Table 3 confirms these observations when taking into account the mean average precision of the ResNet-50 and the convolutional features of networks pre-trained on ImageNet. The difference is more than 10 points gained in mean average precision even when comparing against the ResNet architecture. This means that a ResNet fully optimized for image matching captures the visual information much better than just the convolutional features of a pre-trained network. This is expected, since that was one of the goals of the approach proposed by Gordo et al. [10].

Another observation we can make from Table 3 is that fine-tuning a network on the reference dataset consistently out-performs a pre-trained network. This shows that transfer learning is very powerful for small datasets with many classes. Indeed, networks with many parameters such as AlexNet and ResNet could not have been trained on such small datasets with uninitialized weights.

However, when comparing the classification fine-tuning method with the simplified Siamese architecture (fine-tuning with a triplet loss), it is not as clear which one performs better. From the results, we can see that the classification fine-tuning has a better performance for AlexNet while the triplet loss fine-tuning has a better performance for ResNet-152. This is most likely due to two factors: the hyper-parameters when training the Siamese AlexNet were not perfectly suited, hence the convergence behavior is not as good as with the Siamese ResNet. Furthermore, the AlexNet fine-tuned for classification has a much larger descriptor of dimension 9216 versus the descriptor of dimension 2048 of the simplified Siamese architecture. This may explain that the simplified Siamese architecture performs worse in this case.

Finally, when comparing the proposed architecture with the previous ones, it is clear that the proposed architecture out-performs all of them. It achieves higher precision@1 as well as higher mean average precision, especially when combined with the instance feature augmentation. The comparison with the ResNet-50 from Gordo et al [10] is difficult though. This is because on the one hand, our proposed network is trained on the reference dataset used when comparing images, giving it an unfair advantage. On the other hand, the ResNet-50 is trained on the much larger Landmarks dataset [5], giving it the advantage of data volume. The training methodology developed by Gordo et al. is not applicable to a small, clean dataset, such as the ones used in our evaluation.

The figure 7 shows some example of success and failure of the system. The first two lines are successfully recognize images. The last two are failing examples. Each line represents the query, and the two first images return by the system. The system fail if the first image returned do not represent the same instance than the query. On the two failed example, the system successfully return the correct image, but as second closest image.

# 7.4 Study on Data augmentation

In the previous subsection, we presented the results obtained using all the elements described earlier. We focus here on a detail analysis of the impact of each specific data augmentation, namely rotation, scaling and flipping. Table 4 describes the combinations of these augmentations on the fine-tuned results of the best network tested, namely Resnet-152.

We specifically focus on the impact of the flipping augmentation in our case: as flipping augmentation is commonly used used classical from image classification learning, the flip of images does not seem *a priori* a good idea in our case of instance search, as we do not want to confuse painting that may differ due to flip. As an example, masterpiece "4900 Colors"<sup>1</sup> of Gerhart

<sup>&</sup>lt;sup>1</sup> https://www.gerhard-richter.com/en/art/microsites/4900-colours



Fig. 7 Success and Failing examples. The first column are test queries. The second column are the closest image from the dataset that the system found.

Rotation	Scaling	Flipping	Mean Precision@1 (in %)	
			CLICIDE	GaRoFou
			72.12	92.93
			74.55	94.02
	$\checkmark$		76.36	93.48
			72.12	94.57
$\checkmark$	$\checkmark$		76.97	94.57
		$\checkmark$	75.75	94.57
	$\checkmark$	$\checkmark$	78.18	93.48
	$\checkmark$	$\checkmark$	79.39	94.57

Table 4 Influence of rotation/scaling/flipping data augmentation on ResNet-152 fine-tuning results.

Richter is a good example of such case where the flip may not be considered in the learning set. As table 4 shows, compared to no augmentation (first line of results), rotation-only and scaling-only augmentations increases the learning of the network. However, as expected, the flipping-only augmentation does not increases the quality of the learning. Another interesting finding is that, any pair of combination of augmentations outperform the single component augmentations. Surprisingly, this remarks also holds when considering the flipping augmentation, which is somewhat counterintuitive. The conclusion drawn fro this table is that using all the data augmentations is the best solution, even for paintings as in CLICIDE.drawn fro this table is that using all the data augmentations is the best solution, even for paintings as in CLICIDE. From the experiments on the GaRoFou dataset, we can not conclude about the influence of data augmentation due to the lack of differences between runs. The results are more likely to depend on which local minimum we are.

# **8** Conclusion

This paper presents a novel approach for instance and image retrieval with low variability and small datasets. The proposed approach consists of two key elements. First, we leverage the concept of fully convolutional networks in order to perform classification training at different scales, without a heavy computational overhead. Second, we show that the fully convolutional network can be used to obtain region proposals without the need for an additional component in the network and training. This is particularly important, since region proposals are costly to define manually in our research problem. The region proposals used by the state of the art do not seem applicable to that kind of problem of instance search.

The proposed model consist of first fine-tuning a network over the target dataset. Then a Fully Convolutional Network is trained over the images, each one at different scales. This network is then train in siamese configuration, with a modified triplet loss. This network is used to extract an image representation, which is used to do Instance Retrieval. Finally, the proposed network keeps all the benefits of state-of-the-art approaches: it can be trained end-to-end and it produces an effective global descriptor, which can be compared using the dot product. Additionally, it is modular in the sense that it can be built upon any type of CNN, pre-trained for classification.

Through multiple experiments on two datasets, we show that the descriptor obtained using our proposed network outperforms previous state-of-the-art approaches on the instance search task, while being just as memory-efficient and fast for encoding images. The experiments were conducted on two egocentric image datasets taken from museum visits.

## References

- Arandjelovic, R., Gronat, P., Torii, A., Pajdla, T., Sivic, J.: Netvlad: Cnn architecture for weakly supervised place recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5297–5307 (2016)
- Arandjelović, R., Zisserman, A.: Three things everyone should know to improve object retrieval. In: Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on, pp. 2911–2918. IEEE (2012). URL http://ieeexplore.ieee.org/abstract/document/6248018/
- Babenko, A., Lempitsky, V.: Aggregating local deep features for image retrieval. In: Proceedings of the IEEE international conference on computer vision, pp. 1269–1277 (2015)
- Babenko, A., Slesarev, A., Chigorin, A., Lempitsky, V.: Neural codes for image retrieval. In: European conference on computer vision, pp. 584–599. Springer (2014)
- Babenko, A., Slesarev, A., Chigorin, A., Lempitsky, V.: Neural codes for image retrieval. In: European conference on computer vision, pp. 584–599. Springer (2014)
- Barroso, L.A., Dean, J., Holzle, U.: Web search for a planet: The google cluster architecture. IEEE micro 23(2), 22–28 (2003)
- Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., Shah, R.: Signature verification using a" siamese" time delay neural network. In: Advances in Neural Information Processing Systems, pp. 737–744 (1994)
- 8. Fischer, P., Dosovitskiy, A., Brox, T.: Descriptor matching with convolutional neural networks: a comparison to sift. arXiv preprint arXiv:1405.5769 (2014)
- Gong, Y., Wang, L., Guo, R., Lazebnik, S.: Multi-scale orderless pooling of deep convolutional activation features. In: European conference on computer vision, pp. 392–407. Springer (2014)
- Gordo, A., Almazán, J., Revaud, J., Larlus, D.: Deep Image Retrieval: Learning Global Representations for Image Search. In: Computer Vision – ECCV 2016, pp. 241–257. Springer, Cham (2016)
- 11. Gordo, A., Almazan, J., Revaud, J., Larlus, D.: End-to-end learning of deep visual representations for image retrieval. arXiv preprint arXiv:1610.07940 (2016). URL https://arxiv.org/abs/1610.07940
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016)
- Jegou, H., Douze, M., Schmid, C.: Hamming embedding and weak geometric consistency for large scale image search. Computer Vision–ECCV 2008 pp. 304–317 (2008)
- Krizhevsky, A., Sutskever, I., Hinton, G.E.: ImageNet Classification with Deep Convolutional Neural Networks. In: F. Pereira, C.J.C. Burges, L. Bottou, K.Q. Weinberger (eds.) Advances in Neural Information Processing Systems 25, pp. 1097–1105. Curran Associates, Inc. (2012). URL http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

- Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3431–3440 (2015)
- 16. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. International journal of computer vision **60**(2), 91–110 (2004)
- Paulin, M., Douze, M., Harchaoui, Z., Mairal, J., Perronin, F., Schmid, C.: Local convolutional features with unsupervised training for image retrieval. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 91–99 (2015)
- Perronnin, F., Dance, C.: Fisher kernels on visual vocabularies for image categorization. In: Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on, pp. 1–8. IEEE (2007)
- Philbin, J., Chum, O., Isard, M., Sivic, J., Zisserman, A.: Object retrieval with large vocabularies and fast spatial matching. In: Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on, pp. 1–8. IEEE (2007). URL http://ieeexplore.ieee.org/abstract/document/4270197/
- Philbin, J., Chum, O., Isard, M., Sivic, J., Zisserman, A.: Lost in quantization: Improving particular object retrieval in large scale image databases. In: Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on, pp. 1–8. IEEE (2008). URL http://ieeexplore.ieee.org/abstract/document/4587635/
- Portaz, M., Poignant, J., Budnik, M., Mulhem, P., Chevallet, J., Goeuriot, L.: Construction et évaluation d'un corpus pour la recherche d'instances d'images muséales. In: COnférence en Recherche d'Informations et Applications - CORIA 2017, 14th French Information Retrieval Conference. Marseille, France, March 29-31, 2017. Proceedings, Marseille, France, March 29-31, 2017., pp. 17–34 (2017)
- Radenović, F., Tolias, G., Chum, O.: CNN Image Retrieval Learns from BoW: Unsupervised Fine-Tuning with Hard Examples. In: Computer Vision ECCV 2016, pp. 3–20. Springer, Cham (2016)
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. International Journal of Computer Vision 115(3), 211–252 (2015). DOI 10.1007/s11263-015-0816-y. URL http://link.springer.com/10.1007/s11263-015-0816-y
- Salvador, A., Giro-i Nieto, X., Marques, F., Satoh, S.: Faster R-CNN Features for Instance Search. arXiv:1604.08893 [cs] (2016). URL http://arxiv.org/abs/1604.08893. ArXiv: 1604.08893
- Schroff, F., Kalenichenko, D., Philbin, J.: Facenet: A unified embedding for face recognition and clustering. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 815–823 (2015)
- Sharif Razavian, A., Azizpour, H., Sullivan, J., Carlsson, S.: Cnn features off-the-shelf: an astounding baseline for recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp. 806–813 (2014)
- Simo-Serra, E., Trulls, E., Ferraz, L., Kokkinos, I., Moreno-Noguer, F.: Fracking deep convolutional image descriptors. arXiv preprint arXiv:1412.6537 (2014)
- Sivic, J., Zisserman, A., et al.: Video google: A text retrieval approach to object matching in videos. In: iccv, vol. 2, pp. 1470–1477 (2003)
- Tolias, G., Sicre, R., Jégou, H.: Particular object retrieval with integral max-pooling of CNN activations. arXiv:1511.05879 [cs] (2015). URL http://arxiv.org/abs/1511.05879. ArXiv: 1511.05879
- 30. Turcot, P., Lowe, D.G.: Better matching with fewer features: The selection of useful features in large database recognition problems. In: Computer Vision Workshops (ICCV Workshops), 2009 IEEE 12th International Conference on, pp. 2109–2116. IEEE (2009). URL http://ieeexplore.ieee.org/abstract/document/5457541/
- Yosinski, J., Clune, J., Bengio, Y., Lipson, H.: How transferable are features in deep neural networks? arXiv:1411.1792 [cs] (2014). URL http://arxiv.org/abs/1411.1792. ArXiv: 1411.1792
- Zbontar, J., LeCun, Y.: Stereo matching by training a convolutional neural network to compare image patches. Journal of Machine Learning Research 17(1-32), 2 (2016)