

Deep Feature Factorization For Concept Discovery

Edo Collins¹, Radhakrishna Achanta², and Sabine Süsstrunk¹

¹ School of Computer and Communication Sciences, EPFL

² Swiss Data Science Center, EPFL and ETHZ

{[edo.collins](mailto:edo.collins@epfl.ch),[radhakrishna.achanta](mailto:radhakrishna.achanta@epfl.ch),[sabine.sustrunk](mailto:sabine.sustrunk@epfl.ch)}@epfl.ch

Abstract. We propose Deep Feature Factorization (DFF), a method capable of localizing similar semantic concepts within an image or a set of images. We use DFF to gain insight into a deep convolutional neural network’s learned features, where we detect hierarchical cluster structures in feature space. This is visualized as heat maps, which highlight semantically matching regions across a set of images, revealing what the network ‘perceives’ as similar. DFF can also be used to perform co-segmentation and co-localization, and we report state-of-the-art results on these tasks.

Keywords: Neural network interpretability, Part co-segmentation, Co-segmentation, Co-localization, Non-negative matrix factorization

1 Introduction

As neural networks become ubiquitous, there is an increasing need to understand and interpret their learned representations [25, 27]. In the context of convolutional neural networks (CNNs), methods have been developed to explain predictions and latent activations in terms of heat maps highlighting the image regions which caused them [37, 31].

In this paper, we present Deep Feature Factorization (DFF), which exploits non-negative matrix factorization (NMF) [22] applied to activations of a deep



Fig. 1: *What in this picture is the same as in the other pictures?* Our method, Deep Feature Factorization (DFF), allows us to see how a deep CNN trained for image classification would answer this question. (a) Pyramids, animals and people correspond across images. (b) Monument parts match with each other.

CNN layer to find semantic correspondences across images. These correspondences reflect semantic similarity as indicated by clusters in a deep CNN layer feature space. In this way, we allow the CNN to show us which image regions it ‘thinks’ are similar or related across a set of images as well as within a single image. Given a CNN, our approach to semantic *concept discovery* is unsupervised, requiring only a set of input images to produce correspondences. Unlike previous approaches [2, 11], we do not require annotated data to detect semantic features. We use annotated data for evaluation only.

We show that when using a deep CNN trained to perform ImageNet classification [30], applying DFF allows us to obtain heat maps that correspond to semantic concepts. Specifically, here we use DFF to localize objects or object parts, such as the *head* or *torso* of an animal. We also find that parts form a hierarchy in feature space, e.g., the activations cluster for the concept *body* contains a sub-cluster for *limbs*, which in turn can be broken down to *arms* and *legs*. Interestingly, such meaningful decompositions are also found for object classes never seen before by the CNN.

In addition to giving an insight into the knowledge stored in neural activations, the heat maps produced by DFF can be used to perform co-localization or co-segmentation of objects and object parts. Unlike approaches that delineate the common object across an image set, our method is also able to retrieve distinct parts *within* the common object. Since we use a pre-trained CNN to accomplish this, we refer to our method as performing weakly-supervised co-segmentation.

Our main contribution is introducing Deep Feature Factorization as a method for semantic concept discovery, which can be used both to gain insight into the representations learned by a CNN, as well as to localize objects and object parts within images. We report results on several datasets and CNN architectures, showing the usefulness of our method across a variety of settings.

2 Related work

2.1 Localization with CNN Activations

Methods for the interpretation of hidden activations of deep neural networks, and in particular of CNNs, have recently gained significant interest [25]. Similar to DFF, methods have been proposed to localize objects within an image by means of heat maps [37, 31].

In these works [37, 31], localization is achieved by computing the importance of convolutional feature maps with respect to a particular output unit. These methods can therefore be seen as supervised, since the resulting heat maps are associated with a designated output unit, which corresponds to an object class from a predefined set. With DFF, however, heat maps are *not* associated with an output unit or object class. Instead, DFF heat maps capture common activation patterns in the input, which additionally allows us to localize objects never seen before by the CNN, and for which there is no relevant output unit.

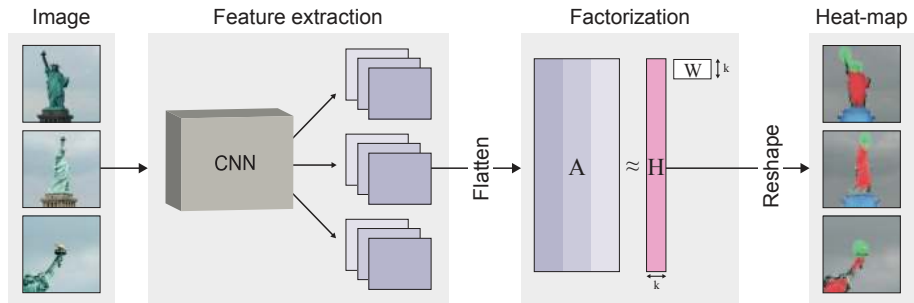


Fig. 2: An illustration of Deep Feature Factorization. We extract features from a deep CNN and view them as a matrix. We apply NMF to the feature matrix and reshape the resulting k factors into k heat maps. See section 3 for a detailed explanation. Shown: Statute of Liberty subset from iCoseg with $k = 3$.

2.2 CNN Features as Part Detectors

The ability of DFF to localize parts stems from the CNN’s ability to distinguish parts in the first place. In Gonzales et al. [11] and Bau et al. [2] the authors attempt to detect learned part-detectors in CNN features, to see if such detectors emerge, even when the CNN is trained with object-level labels. They do this by measuring the overlap between feature map activations and ground truth labels from a part-level segmentation dataset. The availability of ground truth is essential to their analysis, yielding a catalog of CNN units that sufficiently correspond to labels in the dataset.

We confirm their observations that part detectors do indeed emerge in CNNs. However, as opposed to these previous methods, our NMF-based approach does not rely on ground truth labels to find the parts in the input. We use labeled data for evaluation only.

2.3 Non-negative Matrix Factorization

Non-negative matrix factorization (NMF) has been used to analyze data from various domains, such as audio source separation [12], document clustering [36], and face recognition [13].

There has been work extending NMF to multiple layers [6], implementing NMF using neural networks [9] and using NMF approximations as input to a neural network [34]. However, to the best of our knowledge, the application of NMF to the activations of a pre-trained neural network, as is done in DFF, has not been previously proposed.

3 Method

3.1 CNN Feature Space

In the context of CNNs, an input image \mathcal{I} is seen as a tensor of dimension $h_{\mathcal{I}} \times w_{\mathcal{I}} \times c_{\mathcal{I}}$, where the first two dimensions are the height and the width of the image, respectively, and the third dimension is the number of color channels, e.g., 3 for RGB. Viewed this way, the first two dimensions of \mathcal{I} can be seen as a spatial grid, with the last dimension being a $c_{\mathcal{I}}$ -dimensional feature representation of a particular spatial position. For an RGB image, this feature corresponds to color.

As the the image gets processed layer by layer, the hidden activation at the ℓ th layer of the CNN is a tensor we denote $\mathcal{A}_{\mathcal{I}}^{\ell}$ of dimension $h_{\ell} \times w_{\ell} \times c_{\ell}$. Notice that generally $h_{\ell} < h_{\mathcal{I}}$, $w_{\ell} < w_{\mathcal{I}}$ due to pooling operations commonly used in CNN pipelines. The number of channels c_{ℓ} is user-defined as part of the network architecture, and in deep layers is often on the order of 256 or 512.

The tensor $\mathcal{A}_{\mathcal{I}}^{\ell}$ is also called a *feature map* since it has a spatial interpretation similar to that of the original image \mathcal{I} : the first two dimensions represent a spatial grid, where each position corresponds to a *patch* of pixels in \mathcal{I} , and the last dimension forms a c_{ℓ} -dimensional representation of the patch. The intuition behind deep learning suggests that the deeper layer ℓ is, the more abstract and semantically meaningful are the c_{ℓ} -dimensional features [3].

Since a feature map represents multiple patches (depending on the size of image \mathcal{I}), we view them as points inhabiting the same c_{ℓ} -dimensional space, which we refer to as the CNN *feature space*. Having potentially many points in that space, we can apply various methods to find directions that are ‘interesting’.

3.2 Matrix Factorization

Matrix factorization algorithms have been used for data interpretation for decades. For a data matrix A , these methods retrieve an approximation of the form:

$$\begin{aligned} A &\approx \hat{A} = HW & (1) \\ \text{s.t. } A, \hat{A} &\in \mathcal{R}^{n \times m}, H \in \mathcal{R}^{n \times k}, W \in \mathcal{R}^{k \times m} \end{aligned}$$

where \hat{A} is a low-rank matrix of a user-defined rank k . A data point, i.e., a row of A , is explained as a weighted combination of the factors which form the rows of W .

A classical method for dimensionality reduction is principal component analysis (PCA) [18]. PCA finds an optimal k -rank approximation (in the ℓ^2 sense) by solving the following objective:

$$\begin{aligned} \text{PCA}(A, k) = \underset{\hat{A}_k}{\text{argmin}} \quad & \|A - \hat{A}_k\|_F^2, & (2) \\ \text{subject to} \quad & \hat{A}_k = AV_kV_k^{\top}, V_k^{\top}V_k = I_k, \end{aligned}$$

where $\|\cdot\|_F$ denotes the Frobenius norm and $V_k \in \mathcal{R}^{m \times k}$. For the form of Eq. (1), we set $H = AV_k$, $W = V_k^{\top}$. Note that the PCA solution generally contains

negative values, which means the combination of PCA factors (i.e., principal components) leads to the canceling out of positive and negative entries. This cancellation makes intuitive interpretation of individual factors difficult.

On the other hand, when the data A is non-negative, one can perform non-negative matrix factorization (NMF):

$$\begin{aligned} \text{NMF}(A, k) = \operatorname{argmin}_{\hat{A}_k} \quad & \|A - \hat{A}_k\|_F^2, \\ \text{subject to} \quad & \hat{A}_k = HW, \forall ij, H_{ij}, W_{ij} \geq 0, \end{aligned} \quad (3)$$

where $H \in \mathcal{R}^{n \times k}$ and $W \in \mathcal{R}^{k \times m}$ enforce the dimensionality reduction to rank k . Capturing the structure in A while forcing combinations of factors to be additive results in factors that lend themselves to interpretation [22].

3.3 Non-negative Matrix Factorization on CNN Activations

Many modern CNNs make use of the rectified linear activation function, $\max(x, 0)$, due to its desirable gradient properties. An obvious property of this function is that it results in non-negative activations. NMF is thus naturally applicable in this case.

Recall the activation tensor for image \mathcal{I} and layer ℓ :

$$\mathcal{A}_{\mathcal{I}}^{\ell} \in \mathbb{R}^{h \times w \times c} \quad (4)$$

where \mathbb{R} refers to the set of non-negative real numbers. To apply matrix factorization, we partially flatten \mathcal{A} into a matrix whose first dimension is the product of h and w :

$$A_{\mathcal{I}}^{\ell} \in \mathbb{R}^{(h \cdot w) \times c} \quad (5)$$

Note that the matrix $A_{\mathcal{I}}^{\ell}$ is effectively a ‘bag of features’ in the sense that the spatial arrangement has been lost, i.e., the rows of $A_{\mathcal{I}}^{\ell}$ can be permuted without affecting the result of factorization. We can naturally extend factorization to a set of n images, by vertically concatenating their features together:

$$A = \begin{bmatrix} A_1^{\ell} \\ \vdots \\ A_n^{\ell} \end{bmatrix} \in \mathbb{R}^{(n \cdot h \cdot w) \times c} \quad (6)$$

For ease of notation we assumed all images are of equal size, however, there is no such limitation as images in the set may be of any size. By applying NMF to A we obtain the two matrices from Eq. 1, $H \in \mathbb{R}^{(n \cdot h \cdot w) \times k}$ and $W \in \mathbb{R}^{k \times c}$.

3.4 Interpreting NMF Factors

The result returned by the NMF consists of k *factors*, which we will call DFF factors, where k is the predefined rank of the approximation.

The W Matrix Each row W_j ($1 \leq j \leq k$) forms a c -dimensional vector in the CNN feature space. Since NMF can be seen as performing clustering [8], we view a factor W_j as a centroid of an activation cluster, which we show corresponds to coherent object or object-part.

The H Matrix The matrix H has as many rows as the activation matrix A , one corresponding to every spatial position in every image. Each row H_i holds coefficients for the weighted sum of the k factors in W , to best approximate the c -dimensional A_i .

Each column H_j ($1 \leq j \leq k$) can be reshaped into n **heat maps** of dimension $h \times w$, which highlight regions in each image that correspond to the factor W_j . These heat maps have the same spatial dimensions as the CNN layer which produced the activations, often low. To match the size of the heat map with the input image, we upsample it with bilinear interpolation.

4 Experiments

In this section we first show that DFF can produce a hierarchical decomposition into semantic parts, even for sets of very few images (section 4.3). We then move on to larger-scale, realistic datasets where we show that DFF can perform state-of-the-art weakly-supervised object co-localization and co-segmentation, in addition to part co-segmentation (sections 4.4 and 4.5).

4.1 Implementation Details

NMF. NMF optimization with multiplicative updates [23] relies on dense matrix multiplications, and can thus benefit from fast GPU operations. Using an NVIDIA Titan X, our implementation of NMF can process over 6K images of size 224×224 at once with $k = 5$, and requires less than a millisecond per image. Our code is available online.

Neural Network Models. We consider five network architectures in our experiments, namely VGG-16 and VGG-19 [32], with and without batch-normalization [17], as well as ResNet-101 [16]. We use the publicly available models from [26].

4.2 Segmentation and localization methods

In addition to gaining insights into CNN feature space, DFF has utility for various tasks with subtle but important differences in naming:

- **Segmentation vs. Localization** is the difference between predicting pixel-wise binary masks and predicting bounding boxes, respectively.
- **Segmentation vs. co-segmentation** is the distinction between segmenting a single image into regions and jointly segmenting multiple images, thereby producing a correspondence between regions in different images (e.g., *cats* in all images belong to the same segment).

- **Object co-segmentation vs. Part co-segmentation.** Given a set of images representing a common object, the former performs binary background-foreground separation where the foreground segment encompasses the entirety of the common object (e.g., *cat*). The latter, however, produces k segments, each corresponding to a *part* of the common object (e.g., *cat head*, *cat legs*, etc.).

When applying DFF with $k = 1$ can we compare our results against object co-segmentation (background-foreground separation) methods and object co-localization methods.

In section 4.3 we compare DFF against three state-of-the-art co-segmentation methods. The supervised method of Vicente et al. [33] chooses among multiple segmentation proposals per image by learning a regressor to predict, for pairs of images, the overlap between their proposals and the ground truth. Input to the regressor included per-image features, as well as pairwise features. The methods Rubio et al. [29] and Rubinstein et al. [28] are unsupervised and rely on a Markov random field formulation, where the unary features are based on surface image features and various saliency heuristics. For pairwise terms, the former method uses a per-image segmentation into regions, followed by region-matching across images. The latter approach uses a dense pairwise correspondence term between images based on local image gradients.

In section 4.4 we compare against several state-of-the-art object co-localization methods. Most of these methods operate by selecting the best of a set of object proposals, produced by a pre-trained CNN [24] or an object-saliency heuristic [5, 19]. The authors of [21] present a method for unsupervised object co-localization that, like ours, also makes use of CNN activations. Their approach is to apply k -means clustering to globally max-pooled activations, with the intent of clustering all highly active CNN filters together. Their method therefore produces a *single* heat map, which is appropriate for object co-segmentation, but *cannot* be extended to part co-segmentation.

When $k > 1$, we use DFF to perform part co-segmentation. Since we have not come across examples of part co-segmentation in the literature, we compare against a method for supervised part segmentation, namely Wang et al. [35] (Table 3 in section 4.5). Their method relies on a compositional model with strong explicit priors w.r.t to part size, hierarchy and symmetry. We also show results for two baseline methods described in [35]: PartBB+ObjSeg where segmentation masks are produced by intersecting part-bounding-boxes [4] with whole-object segmentation masks [14]. The method PartMask+ObjSeg is similar, but here bounding-boxes are replaced with the best of 10 pre-learned part masks.

4.3 Experiments on iCoseg

Dataset The iCoseg dataset [1] is a popular benchmark for co-segmentation methods. As such, it consists of 38 sets of images, where each image is annotated with a pixel-wise mask encompassing the main object common to the set. Images within a set are uniform in that they were all taken on a single occasion, depicting

the same objects. The challenging aspect of this datasets lies in the significant variability with respect to viewpoint, illumination, and object deformation.

We chose five sets and further labeled them with pixel-wise object-part masks (see Table 1). This process involved partitioning the given ground truth mask into sub-parts. We also annotated common background objects, e.g., *camel* in the *Pyramids* set (see Figure 1). Our part-annotation for iCoseg is available online. The number of images in these sets ranges from as few as 5 up to 41. When comparing against [33] and [29] in Table 1, we used the subset of iCoseg used in those papers.

Part co-segmentation For each set in iCoseg, we obtained activations from the deepest convolutional layer of VGG19 (`conv5_4`), and applied NMF to these activations with increasing values of k . The resulting heat maps can be seen in Figures 1 and 3.

Qualitatively, we see a clear correspondence between DFF factors and coherent object-parts, however, the heat maps are coarse. Due to the low resolution of deep CNN activations, and hence of the heat map, we get blobs that do not perfectly align with the underlying region of interest. We therefore also report additional results with a post-processing step to refine the heat maps, described below.

We notice that when $k = 1$, the single DFF factor corresponds to a whole object, encompassing multiple object-parts. This, however, is not guaranteed, since it is possible that for a set of images, setting $k = 1$ will highlight the *background* rather than the foreground. Nonetheless, as we increase k , we get a decomposition of the object or scene into individual parts. This behavior reveals a hierarchical structure in the clusters formed in CNN feature space.

For instance, in Figure 3 (a), we can see that $k = 1$ encompasses most of gymnast’s body, $k = 2$ distinguished her midsection from her limbs, $k = 3$ adds a finer distinctions between arms and legs, and finally $k = 4$ adds a new component that localizes the beam. This observation also indicates the CNN has learned representation that ‘explains’ these concepts with invariance to pose, e.g., leg positions in the 2nd, 3rd, and 4th columns.

A similar decomposition into legs, torso, back, and head can be seen for the elephants in Figure 3 (b). This shows that we can localize different objects and parts even when they are all common across the image set. Interestingly, the decompositions shown in Figure 1 exhibit similar high semantic quality in spite of their dissimilarity to the ImageNet training data, as neither pyramids nor the Taj Mahal are included as class labels in that dataset. We also note that as some of the given sets contain as few as 5 images (Figure 1 (b) comprises the whole set), our method does not require many images to find meaningful structure.

Object and Part co-segmentation We operationalize DFF to perform co-segmentation. To do so we have to first annotate the factors as corresponding to specific ground-truth parts. This can be done manually (as in Table 3) or

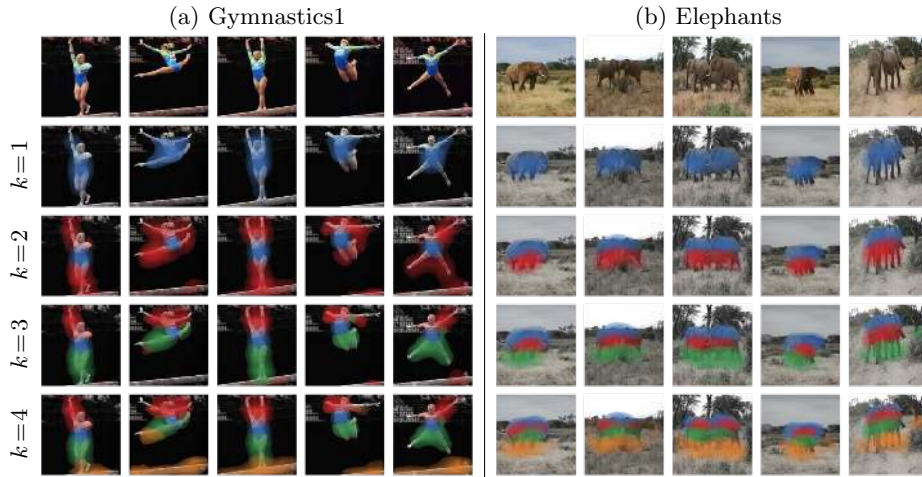


Fig. 3: Example DFF heat maps for images of two sets from iCoseg. Each row shows a separate factorization where the number of DFF factors k is incremented. Different colors correspond to the heat maps of the k different factors. DFF factors correspond well to distinct object parts. This Figure visualizes the data in Table 1, where heat map color corresponds with row color. (Best viewed electronically with a color display)

automatically given ground truth, as described below. We report the intersection-over-union (IoU) score of each factor with its associated parts in Table 1.

Since the heat maps are of low-resolution, we refine them with post processing. We define a dense conditional random field (CRF) over the heat maps. We use the filter-based mean field approximate inference [20], where we employ guided filtering [15] for the pairwise term, and use the bilinearly upsampled DFF heat maps as unary terms. We refer to DFF with post-processing 'DFF-CRF'.

Each heat map is converted to a binary mask using a thresholding procedure. For a specific DFF factor f ($1 \leq f \leq k$), let $\{H(f, 1), \dots, H(f, n)\}$ be the set of n heat maps associated with n input images, The value of a pixel in the binary map $B(f, i)$ of factor f and image i is 0 if its intensity is lower than the 75th percentile of entries in the set of heat maps $\{H(f, j) | 1 \leq j \leq n\}$.

We associate parts with factors by considering how well a part is covered by a factor's binary masks. We define the *coverage* of part p by factor f as:

$$Cov_{f,p} = \frac{|\sum_i B(f, i) \cap P(p, i)|}{|\sum_i P(p, i)|} \quad (7)$$

The coverage is the percentage of pixels belonging to p that are set to 1 in the binary maps $\{B(f, i) | 1 \leq i \leq n\}$. We associate the part p with factor f when $Cov_{f,p} > Cov_{th}$. We experimentally set the threshold $Cov_{th} = 0.5$.

Finally, we measure the IoU between a DFF factor f and its m associated ground-truth parts $\{p_1^{(f)}, \dots, p_m^{(f)}\}$ similarly to [2], specifically by considering

Method	Elephants	Taj Mahal	Pyramids	Gymnastics1	Statue of Liberty
Object co-segmentation					
Vicente [33]	<i>whole</i> 43	<i>whole</i> 91	-	-	<i>whole</i> 94
Rubio [29]	<i>whole</i> 75	<i>whole</i> 89	-	-	<i>whole</i> 92
Rubinstein [28]	<i>whole</i> 63	<i>whole</i> 48	<i>whole</i> 57	<i>whole</i> 94	<i>whole</i> 70
DFF, $k=1$	<i>whole</i> 65	<i>whole</i> 41	<i>whole</i> 57	<i>whole</i> 43	<i>whole</i> 49
DFF-CRF, $k=1$	<i>whole</i> 76	<i>whole</i> 51	<i>whole</i> 70	<i>whole</i> 52	<i>whole</i> 62
Part co-segmentation					
DFF, $k=2$	<i>torso/back/head</i> 59	<i>dome</i> 33	<i>animal</i> 36	<i>torso/waist</i> 35	<i>torso</i> 36
	<i>torso/leg</i> 35	<i>tower/building</i> 46	<i>pyramid</i> 56	<i>arm/leg/head</i> 20	<i>torch/base/head</i> 28
DFF, $k=3$	<i>back/head</i> 46	<i>building</i> 45	<i>background</i> 27	<i>torso/waist</i> 38	<i>base</i> 14
	<i>torso</i> 25	<i>dome</i> 40	<i>pyramid</i> 55	<i>arm/head</i> 22	<i>torso</i> 39
	<i>leg</i> 21	<i>tower</i> 13	<i>animal</i> 36	<i>leg</i> 33	<i>torch/head</i> 23
DFF, $k=4$	<i>torso/back/head</i> 58	<i>building</i> 72	<i>background</i> 27	<i>torso/waist</i> 40	<i>torso</i> 39
	<i>head</i> 36	<i>dome</i> 43	<i>pyramid</i> 52	<i>torso/arm/head</i> 33	<i>background</i> 44
	<i>torso</i> 20	<i>background</i> 08	<i>animal</i> 37	<i>leg</i> 37	<i>torch/head</i> 26
	<i>leg</i> 16	<i>tower</i> 16	<i>person</i> 12	<i>background</i> 14	<i>base</i> 40

Table 1: Object and part discovery and segmentation on five iCoseg image sets. Part-labels are automatically assigned to DFF factors, and are shown with their corresponding *IoU*-scores. Our results show that clusters in CNN feature space correspond to coherent parts. More so, they indicate the presence of a cluster hierarchy in CNN feature space, where part-clusters can be seen as sub-clusters within object-clusters (See Figures 1, 2 and 3 for visual comparison. Row color corresponds with heat map color). With $k = 1$, DFF can be used to perform object co-segmentation, which we compare against state-of-the-art methods. With $k > 1$ DFF can be used to perform part co-segmentation, which current co-segmentation methods are not able to do.

the dataset-wide *IoU* :

$$P_f(i) = \bigcup_j^m P_j^{(f)} \quad (8)$$

$$IoU_{f,p} = \frac{|\sum_i B_i \cap P_f(i)|}{|\sum_i B_i \cup P_f(i)|} \quad (9)$$

In the top of Table 1 we report results for object co-segmentation ($k = 1$) and show that our method is comparable with the supervised approach of [33] and domain-specific methods of [29] and [28].

The bottom of Table 1 shows the labels and *IoU*-scores for part co-segmentation on the five image sets of iCoseg that we have annotated. These scores correspond to the visualizations of Figures 1 and 3 and confirm what we observe qualitatively.

We can characterize the quality of a factorization as the average *IoU* of each factor with its single best matching part (which is not the background). In Figure 4 (a) we show the average *IoU* for different layer of VGG-19 on iCoseg as the value of k increases. The variance shown is due to repeated trials with

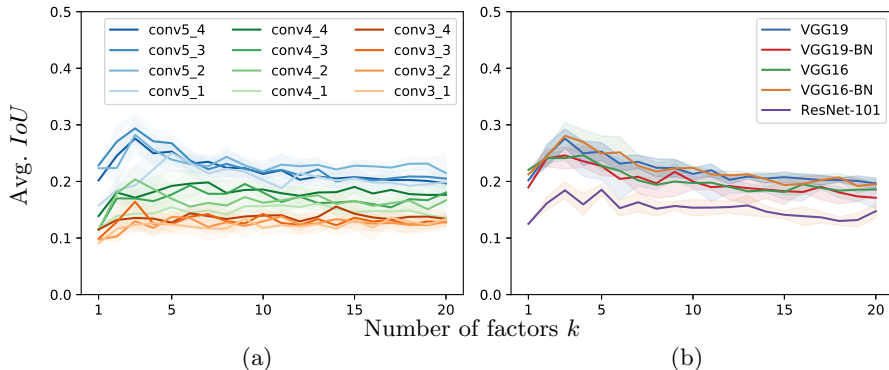


Fig. 4: (a) Average IoU score for DFF on iCoseg. for (a) different VGG19 layers and (b) the deepest convolutional layer for other CNN architectures. Expectedly, different convolutional blocks show a clear difference in matching up with semantic parts, as CNN features capture more semantic concepts. The optimal value for k is data dependent but is usually below 5. We see also that DFF performance is relatively uniform for the VGG family of models.

different NMF initializations. There is a clear gap between convolutional blocks. Performance within a block does not strictly follow the linear order of layers.

We also see that the optimal value for k is between 3 and 5. While this naturally varies for different networks, layers, and data batches, another deciding factor is the resolution of the part ground truth. As k increases, DFF heat maps become more localized, highlighting regions that are beyond the granularity of the ground truth annotation, e.g., a pair of factors that separates *leg* into *ankle* and *thigh*. In Figure 4 (b) we show that DFF performs similarly within the VGG family of models. For ResNet-101 however, the average IoU is distinctly lower.

4.4 Object Co-Localization on PASCAL VOC 2007

Dataset PASCAL VOC 2007 has been commonly used to evaluate whole object co-localization methods. Images in this dataset often comprise several objects of multiple classes from various viewpoints, making it a challenging benchmark. As in previous work [21, 5, 19], we use the *trainval* set for evaluation and filter out images that only contain objects which are marked as *difficult* or *truncated*. The final set has 20 image sets (one per class), with 69 to 2008 images each.

Evaluation The task of co-localization involves fitting a bounding box around the common object in a set of image. With $k = 1$, we expect DFF to retrieve a heat map which localizes that object.

As described in the previous section, after optionally filtering DFF heat maps using a CRF, we convert the heat maps to binary segmentation masks. We follow

[31] and extract a single bounding box per heat map by fitting a box around the largest connected component in the binary map.

We report the standard CorLoc score [7] of our localization. The CorLoc score is defined as the percentage of predicted bounding boxes for which there exists a matching ground truth bounding box. Two bounding boxes are deemed matching if their *IoU* score exceeds 0.5.

The results of our method are shown in Table 2, along with previous methods (described in section 4.2). Our method compares favorably against previous approaches. For instance, we improve co-localization for the class *dog* by 16% higher CorLoc and achieve better co-localization on average, in spite of our approach being simpler and more general.

Method	aero	bicy	bird	boa	bot	bus	car	cat	cha	cow	dtab	dog	hors	mbik	pers	plnt	she	sofa	tra1	tv	mean
Joulin [19]	338	17	21	18	5	27	33	41	6	29	35	32	26	40	18	12	25	28	36	12	25.60
Cho [5]	50	43	30	19	4	62	65	43	9	49	12	44	64	57	15	9	31	34	62	32	36.60
Li [24]	73	45	43	28	7	53	58	45	6	48	14	47	69	67	24	13	52	26	65	17	40.00
Le (A) [21]	70	52	44	30	5	56	60	59	6	49	16	51	59	67	23	12	47	27	59	16	40.36
Le (V) [21]	72	62	48	28	12	64	59	72	6	37	12	45	67	72	19	11	37	29	67	23	41.97
DFF	61	49	54	20	10	60	46	79	4	51	32	67	66	70	19	15	40	32	66	20	42.87
DFF-CRF	64	47	50	16	10	62	52	75	8	53	35	65	65	72	16	14	41	36	63	30	43.51

Table 2: Co-localization results for PASCAL VOC 2007 with DFF $k = 1$. Numbers indicate CorLoc scores. Overall, we exceed the state-of-the-art approaches using a much simpler method.

4.5 Part Co-segmentation in PASCAL-Parts

Dataset The PASCAL-Part dataset [4] is an extension of PASCAL VOC 2010 [10] which has been further annotated with part-level segmentation masks and bounding boxes. The dataset decomposes 16 object classes into fine grained parts, such as *bird-beak* and *bird-tail* etc. After filtering out images containing objects marked as *difficult* and *truncated*, the final set consists of 16 image sets with 104 to 675 images each.

Evaluation In Table 3 we report results for the two classes, *cow* and *horse*, which are also part-segmented by Want et al. as described in section 4.2. Since their method relies on strong explicit priors w.r.t to part size, hierarchy, and symmetry, and its explicit objective is to perform part-segmentation, their results serve as an upper bound to ours. Nonetheless we compare favorably to their results and even surpass them in one case, despite our method not using any hand-crafted features or supervised training.

For this experiment, our strategy for mapping DFF factors ($k = 3$) to their appropriate part labels was with semi-automatic labeling, i.e., we qualitatively examined the heat maps of *only five images*, out of approximately 140 images, and labeled factors as corresponding to the labels shown in Table 3.

In Table 4 we give *IoU* results for five additional classes from PASCAL-Parts, which have been automatically mapped to parts as in section 4.3. In Figure 5

Method	<i>cow</i>			<i>horse</i>		
	<i>head</i>	<i>neck+torso</i>	<i>leg</i>	<i>head</i>	<i>neck+torso</i>	<i>leg</i>
PartBB+ObjSeg	26.77	53.79	11.18	37.32	60.35	27.47
PartMask+ObjSeg	33.19	56.69	11.31	41.84	63.31	21.38
Compositional model [35]	41.55	60.98	30.98	47.21	66.74	38.18
DFF	40.53	59.48	21.57	28.85	54.77	28.94
DFF-CRF	45.20	58.87	24.60	31.05	53.18	28.81

Table 3: Avg. IoU(%) for three fully supervised methods reported in [35] (see section 4.2 for details) and for our weakly-supervised DFF approach. As opposed to DFF, previous approaches shown are fully supervised. Despite not using hand-crafted features, DFF compares favorably to these approaches, and is not specific to these two image classes. We semi-automatically mapped DFF factors ($k = 3$) to their appropriate part labels by examining the heat maps of *only five* images, out of approximately 140 images. This illustrates the usefulness of DFF co-segmentation for fast semi-automatic labeling. See visualization for *cow* heat maps in Figure 5.

k	aeroplane	bird	car	motorbike	cat
1	42 <i>bird</i>	40 <i>car</i>	29 <i>wheel</i>	30 <i>eye/head/neck/nose</i>	31
2	2 <i>beak/eye/head/neck</i>	13 <i>wheel</i>	10 <i>wheel</i>	38 <i>torso</i>	24
3	49 <i>neck/torso/wing</i>	39 <i>door/roof/window</i>	22 <i>person</i>	9 <i>eye/head/neck/nose</i>	36
4	2 <i>leg</i>	2 <i>wheel</i>	10 <i>wheel</i>	30 <i>eye/head/neck/nose</i>	32
5	47 <i>neck/torso/wing</i>	43 <i>door/headlight/licenseplate</i>	24 <i>headlight</i>	1 <i>torso</i>	30
6	35 <i>beak/eye/head/neck/torso</i>	30 <i>mirror/roof/window</i>	20 <i>wheel</i>	29 <i>car/eye/head/neck/nose</i>	38
7	1 <i>foot/leg</i>	3 <i>wheel</i>	9 <i>wheel</i>	33 <i>eye/head/nose</i>	31
8	44 <i>neck/torso/wing</i>	44 <i>headlight/licenseplate</i>	31 <i>person</i>	10 <i>eye/neck/nose</i>	5
9	21 <i>beak/eye/head/neck/torso</i>	30 <i>front</i>	8 <i>wheel</i>	17 <i>car/eye/head/nose</i>	35
10	32 <i>neck</i>	2 <i>mirror/roof/window</i>	22 <i>background</i>	13 <i>torso</i>	27

Table 4: IoU of DFF heat maps with PASCAL-Parts segmentation masks. Each DFF factor is automatically labeled with part labels as in section 4.3. Higher values of k allow DFF to localize finer regions across the image set, some of which go beyond the resolution of the ground truth part annotation. Figure 5 visualizes the results for $k = 3$ (row color corresponds to heat map color).

we visualize these DFF heat maps for $k = 3$, as well as for *cow* from Table 3. When comparing the heat maps against their corresponding *IoU*-scores, several interesting conclusions can be made. For instance, in the case of *motorbike*, the first and third factors for $k = 3$ in Table 4 both seems to correspond with wheel. The visualization in Figure 5 (e) reveals that these factors in fact sub-segment the wheel into top and bottom, which is beyond the resolution of the ground truth data.

We can see also that while the first factor of the class *aeroplane* (Figure 5 (a)) consistently localizes airplane wheels, it does not to achieve high *IoU* due to the coarseness of the heat map.

Returning to Table 4, when $k = 4$, a factor emerges that localizes instances of the class *person*, which occur in 60% of motorbike images. This again shows that while most co-localization methods only describe objects that are common across the image set, our DFF approach is able to find distinctions *within* the set of common objects.

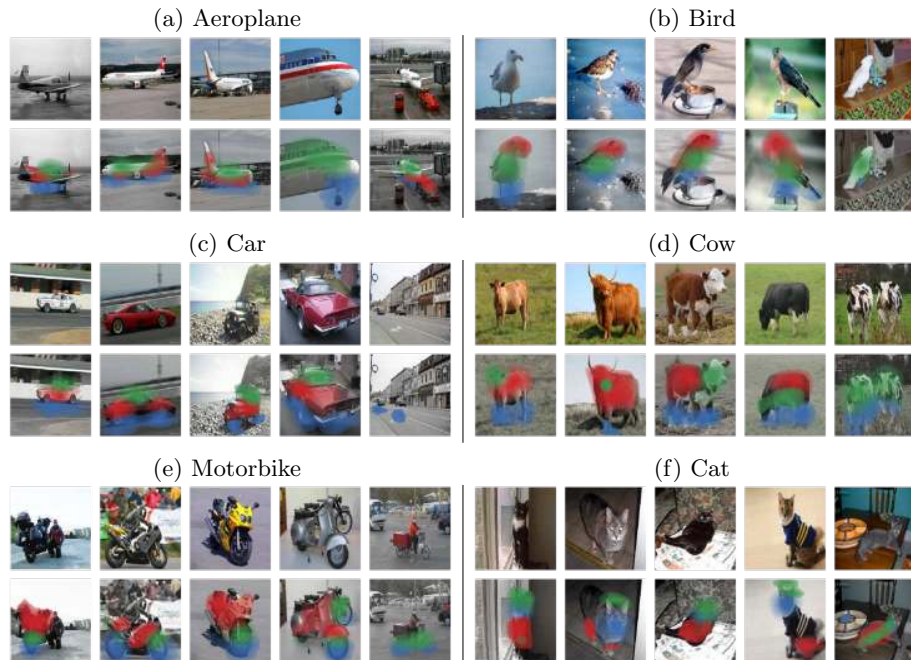


Fig. 5: Example DFF heat maps for images of six classes from PASCAL-Parts with $k = 3$. For each class we show four images that were successfully decomposed into parts, and a failure case on the right. DFF manages to retrieve interpretable decompositions in spite of the great variation in the data. In addition to the DFF factors for *cow* from Table 3, here visualized are the factors which appear in Table 4, where heat map colors correspond to row colors.

5 Conclusions

In this paper, we have presented Deep Feature Factorization (DFF), a method that is able to locate semantic concepts in individual images and across image sets. We have shown that DFF can reveal interesting structures in CNN feature space, such as hierarchical clusters which correspond to a part-based decomposition at various levels of granularity.

We have also shown that DFF is useful for co-segmentation and co-localization, achieving results on challenging benchmarks which are on par with state-of-the-art methods, and can be used to perform semi-automatic image labeling. Unlike previous approaches, DFF can also perform *part* co-segmentation as well, making fine distinction *within* the common object, e.g. matching *head* to *head* and *torso* to *torso*.

References

1. Batra, D., Kowdle, A., Parikh, D., Luo, J., Chen, T.: icoseg: Interactive co-segmentation with intelligent scribble guidance. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 3169–3176. IEEE (2010)
2. Bau, D., Zhou, B., Khosla, A., Oliva, A., Torralba, A.: Network dissection: Quantifying interpretability of deep visual representations. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 3319–3327. IEEE (2017)
3. Bengio, Y., Courville, A., Vincent, P.: Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **35**(8), 1798–1828 (2013)
4. Chen, X., Mottaghi, R., Liu, X., Fidler, S., Urtasun, R., Yuille, A.: Detect what you can: Detecting and representing objects using holistic models and body parts. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 1971–1978 (2014)
5. Cho, M., Kwak, S., Schmid, C., Ponce, J.: Unsupervised object discovery and localization in the wild: Part-based matching with bottom-up region proposals. In: *Computer Vision and Pattern Recognition (CVPR)* (2015)
6. Cichocki, A., Zdunek, R.: Multilayer nonnegative matrix factorisation. *Electronics Letters* **42**(16), 1 (2006)
7. Deselaers, T., Alexe, B., Ferrari, V.: Weakly supervised localization and learning with generic knowledge. *International Journal of Computer Vision (IJCV)* **100**(3), 275–293 (2012)
8. Ding, C., He, X., Simon, H.D.: On the equivalence of nonnegative matrix factorization and spectral clustering. In: *Proceedings of the 2005 SIAM International Conference on Data Mining*. pp. 606–610. SIAM (2005)
9. Dziugaite, G.K., Roy, D.M.: Neural network matrix factorization. arXiv preprint arXiv:1511.06443 (2015)
10. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results. <http://www.pascal-network.org/challenges/VOC/voc2010/workshop/index.html>
11. Gonzalez-Garcia, A., Modolo, D., Ferrari, V.: Do semantic parts emerge in convolutional neural networks? *International Journal of Computer Vision (IJCV)* pp. 1–19 (2017)
12. Grais, E.M., Erdogan, H.: Single channel speech music separation using nonnegative matrix factorization and spectral masks. In: *Digital Signal Processing (DSP)*. pp. 1–6. IEEE (2011)
13. Guillaumet, D., Vitria, J.: Non-negative matrix factorization for face recognition. In: *Topics in artificial intelligence*, pp. 336–344. Springer (2002)
14. Hariharan, B., Arbeláez, P., Girshick, R., Malik, J.: Simultaneous detection and segmentation. *European Conference on Computer Vision (ECCV)* (2014)
15. He, K., Sun, J., Tang, X.: Guided image filtering. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **35**(6), 1397–1409 (2013)
16. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 770–778 (2016)
17. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: *International Conference on Machine Learning (ICML)*. pp. 448–456 (2015)
18. Jolliffe, I.T.: Principal component analysis and factor analysis. In: *Principal component analysis*, pp. 115–128. Springer (1986)

19. Joulin, A., Tang, K., Fei-Fei, L.: Efficient image and video co-localization with frank-wolfe algorithm. In: European Conference on Computer Vision (ECCV). pp. 253–268. Springer (2014)
20. Krähenbühl, P., Koltun, V.: Efficient inference in fully connected crfs with gaussian edge potentials. In: Advances in Neural Information Processing Systems (NIPS). pp. 109–117 (2011)
21. Le, H., Yu, C.P., Zelinsky, G., Samaras, D.: Co-localization with category-consistent features and geodesic distance propagation. In: Computer Vision and Pattern Recognition (CVPR). pp. 1103–1112 (2017)
22. Lee, D.D., Seung, H.S.: Learning the parts of objects by non-negative matrix factorization. *Nature* **401**(6755), 788 (1999)
23. Lee, D.D., Seung, H.S.: Algorithms for non-negative matrix factorization. In: Advances in neural information processing systems. pp. 556–562 (2001)
24. Li, Y., Liu, L., Shen, C., van den Hengel, A.: Image co-localization by mimicking a good detectors confidence score distribution. In: European Conference on Computer Vision (ECCV). pp. 19–34. Springer (2016)
25. Montavon, G., Samek, W., Müller, K.: Methods for interpreting and understanding deep neural networks. *Digital Signal Processing* **73**, 1–15 (2018). <https://doi.org/10.1016/j.dsp.2017.10.011>, <https://doi.org/10.1016/j.dsp.2017.10.011>
26. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)
27. Ribeiro, M.T., Singh, S., Guestrin, C.: Why should I trust you?: Explaining the predictions of any classifier. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining pp. 1135–1144 (2016)
28. Rubinstein, M., Joulin, A., Kopf, J., Liu, C.: Unsupervised joint object discovery and segmentation in internet images. *Computer Vision and Pattern Recognition (CVPR)* (June 2013)
29. Rubio, J.C., Serrat, J., López, A., Paragios, N.: Unsupervised co-segmentation through region matching. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 749–756. IEEE (2012)
30. 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 (IJCV)* **115**(3), 211–252 (2015). <https://doi.org/10.1007/s11263-015-0816-y>
31. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. See <https://arxiv.org/abs/1610.02391> v3 **7**(8) (2016)
32. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
33. Vicente, S., Rother, C., Kolmogorov, V.: Object cosegmentation. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 2217–2224. IEEE (2011)
34. Vu, T.T., Bigot, B., Chng, E.S.: Combining non-negative matrix factorization and deep neural networks for speech enhancement and automatic speech recognition. In: *Acoustics, Speech and Signal Processing (ICASSP)*. pp. 499–503. IEEE (2016)
35. Wang, J., Yuille, A.L.: Semantic part segmentation using compositional model combining shape and appearance. *CVPR* (2015)
36. Xu, W., Liu, X., Gong, Y.: Document clustering based on non-negative matrix factorization. In: Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval. pp. 267–273. ACM (2003)

37. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: Computer Vision and Pattern Recognition (CVPR). pp. 2921–2929. IEEE (2016)