# Analyzing Classifiers: Fisher Vectors and Deep Neural Networks

Sebastian Bach
Fraunhofer HHI
Einsteinufer 37, 10587 Berlin
sebastian.bach@hhi.fraunhofer.de

Alexander Binder
SUTD
8 Somapah Rd, 487372 Singapore
alexander_binder@sutd.edu.sg

Grégoire Montavon
TU Berlin
Marchstr. 23, 10587 Berlin. Germany
gregoire.montavon@tu-berlin.de

Klaus-Robert Müller
TU Berlin
Marchstr. 23, 10587 Berlin. Germany
klaus-robert.mueller@tu-berlin.de

Wojciech Samek
Fraunhofer HHI
Einsteinufer 37, 10587 Berlin
wojciech.samek@hhi.fraunhofer.de

## 1 Abstract

Fisher Vector classifiers and Deep Neural Networks (DNNs) are popular and successful algorithms for solving image classification problems. However, both are generally considered 'black box' predictors as the non-linear transformations involved have so far prevented transparent and interpretable reasoning. Recently, a principled technique, Layer-wise Relevance Propagation (LRP), has been developed in order to better comprehend the inherent structured reasoning of complex nonlinear classification models such as Bag of Feature models or DNNs. In this paper we (1) extend the LRP framework also for Fisher Vector classifiers and then use it as analysis tool to (2) quantify the importance of context for classification, (3) qualitatively compare DNNs against FV classifiers in terms of important image regions and (4) detect potential flaws and biases in data. All experiments are performed on the PASCAL VOC 2007 data set.

## 2 Introduction

Deep neural networks have defined state of the art in many fields, such as image classification [13], image detection [7] and machine translation [25]. While much of research is devoted to extending the applicability of deep neural nets to more domains [8, 12, 30, 11, 10], we focus here on a different question, namely the impact of context, and the ability to use context. This question was raised already during times of the Pascal VOC challenge, where the amount of context was a matter of speculation, c.f. PASCAL VOC workshop presentation slides in [6].

The question of context is considered for two prominent types of classifiers. The first type, Fisher Vectors (FV) [23] are based on computing a single feature map on an image as a whole and subsequently computing one score. In such a setup one can expect that context plays naturally a role for the prediction as the image is processed as a whole during training and test time. In case of small training sample sizes and the absence of opportunities for fine-tuning, Fisher vectors still might be a viable alternative to deep neural nets due to their reduced parameter space. Examples for performance issues of deep neural networks on small sample sizes without finetuning can be seen in [29]. The question of context is also open for the second type, Deep Neural Networks (DNN). One might assume that context plays no role for neural networks when they are used in classification by detection setups. For example, a recent Imagenet challenge winner relied on 144 crops per test image and classifier [26]. Another work using Pascal VOC data [19] used at test time 500 multi-scale patches per test image. However

in certain setups computing several hundred windows as required for classification by detection setups may not be possible, e.g. when using hardware without GPUs and much main memory, such as used consumer laptops or smartphones, and when having time constraints for computation of the test prediction on an image. One can expect to see a larger impact of context when resorting to a few regions of an image at test time only, and thus training and testing with larger image patches.

Our contribution here is as follows. (1) We extend the method of [1] to Fisher vectors, and apply relevance propagation for the first time to Fisher vectors. (2) We define measures for the amount of context used for prediction in a single test image. (3) We apply the measures of context for neural networks and Fisher vector based classifiers on the Pascal VOC dataset, as it offers a way to approximately validate context by its bounding box annotation. We compare the context dependence of Fisher vectors against neural nets which were trained on larger patches of input images. (4) We show that this methodology is able to identify strong cases of context and biases in the training data even without using bounding box information.

The next section reviews related work. Section 4 briefly describes the Fisher Vector classifier. Section 5 introduces the extended LRP method to decompose a Fisher Vector prediction into scores for small regions of the order of a local feature. The same section also proposes a novel LRP-based measure of the importance of context. Section 6 introduces the experimental setup and presents results. The paper concludes in Section 7 with a summary and an outlook.

## 3   Related Work

In recent years, interest in understanding image representations [17, 15, 18] and being able to explain the decision process of a classification system has increased, with e.g., gradient-based sensitivity analysis [2, 24]. However, many approaches have been conceived with a specific pipeline architecture in mind. So do [27] explain predictions for bag of word features with hard mapping (Vector Quantization) and Histogram Intersection kernels, and [16] identifies image regions critical for the prediction of a linear SVM classifier with max-pooling feature aggregation algorithm. A solution especially dedicated to visualize image regions triggering the prediction of deep convolutional neural networks with max-pooling layers and has been proposed in [29]

with *deconvolution nets*.

Recently, a paradigm called Layer-wise Relevance Propagation (LRP) has been introduced in [1] as a way to compute partial prediction contributions – or *relevance values $R$* – for intermediate and input representations based on the final classifier output. It computes scores for regions or pixels of an image explaining the prediction itself rather than the effect of single neurons or particular layers. It is applied in [1] to Bag of visual words classifiers and deep neural networks; in this paper we extend this method to make it applicable to Fisher Vector classifiers.

## 4   Fisher Vectors in a Nutshell

Fisher Vectors [20, 23] are a powerful tool to compute rich image or video representations and provide state-of-the-art performance amongst feature extraction algorithms. Figure 1 summarizes the steps involved in computing FV representation of an image. We introduce here a notation which later will be used in the Section 5.

An integral part for computing FVs is to fit a Gaussian Mixture Model (GMM) on top of the local descriptors $L = \{l\}$ extracted from the training data to serve as a soft vocabulary of visual prototypes. Assuming a $K$-component GMM $\lambda = \{(\pi_k, \mu_k, \Sigma_k)\}_{k=1..K}$, then $\pi_k$ is the mixture weight of component $k$, with $\sum_k \pi_k = 1$ and $\forall k : \pi_k \geq 0$, $\mu_k$ is the mean vector of the $k$th mixture component and $\Sigma_k$ its (diagonal) covariance matrix. For the computation of a *full* FV representation of an image, each local descriptor $l$ is related to all $K$ components of the trained GMM in its 0th (soft mapping weight), 1st (deviation from mean) and 2nd moment (variance) [23]:

$$\Psi_{\pi_k}(l) = \frac{1}{\sqrt{\pi_k}} \left( \gamma_k(l) - \pi_k \right) \qquad (1)$$

$$\Psi_{\mu_k}(l) = \frac{1}{\sqrt{\pi_k}} \gamma_k(l) \left( \frac{l - \mu_k}{\sigma_k} \right) \qquad (2)$$

$$\Psi_{\sigma_k}(l) = \frac{1}{\sqrt{\pi_k}} \gamma_k(l) \frac{1}{\sqrt{2}} \left( \frac{(l - \mu_k)^2}{\sigma_k^2} - 1 \right) \quad (3)$$

with $\Psi_{\pi_k}(l) \in \mathbb{R}$ , both $\Psi_{\mu_k}(l)$ and $\Psi_{\sigma_k}(l) \in \mathbb{R}^D$ and $\gamma_k(l)$ returning the soft assignment of $l$ to the $k$th mixture component. The FV embedding $\Psi_\lambda(l)$ for a single descriptor $l$ is then achieved by concatenating the mapping outputs relative to all $K$ components into a $(1 + 2D)K$ dimensional vector

$$\Psi_\lambda(l) = [\Psi_{\pi_1}(l) \ldots \Psi_{\mu_1}(l) \ldots \Psi_{\sigma_1}(l) \ldots] \quad (4)$$

**Layer 1**
Image of 'bicycle'

Heatmap

**Layer 2**
Local Features

**Layer 3**
GMM fitting          Fisher Vector

$\Psi_\lambda(l) = \left[\;\Psi_{\pi_k}(l)\;\;\Psi_{\mu_k}(l)\;\;\Psi_{\sigma_k}(l)\;\right]$

$\Psi_{\pi_k}(l)\;\Psi_{\mu_k}(l)\;\Psi_{\sigma_k}(l)$

**Layer 4**
Normalization + Linear SVM
(Hellinger's kernel SVM)

$\mathbf{x} = \frac{1}{|L|}\sum_{l \in L}\Psi_\lambda(l)$

$\mathbf{x} \leftarrow sign(\mathbf{x})|\mathbf{x}|^{\frac{1}{2}}$

$\mathbf{x} \leftarrow \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$ $\quad f(\mathbf{x}) = b + \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x})$

Relevance Conservation

$\sum_i R_i^{(1)} = \sum_j R_j^{(2)} \qquad \sum_i R_i^{(2)} = \sum_j R_j^{(3)} \qquad \sum_i R_i^{(3)} = f(\mathbf{x})$

Redistribution Formula

$R_p^{(1)} = \sum_{l \in L(p)} \frac{R_l^{(2)}}{|\text{area}(l)|}$ $\qquad R_l^{(2)} = \sum_d R_d^{(3)} \frac{z_{ld}}{\sum_{l'} z_{l'd} + \epsilon \cdot \text{sign}(\sum_{l'} z_{l'd})}$ $\qquad R_d^{(3)} = \sum_i \alpha_i y_i \phi(x_i)_d \phi(x)_d + \frac{b}{D}$
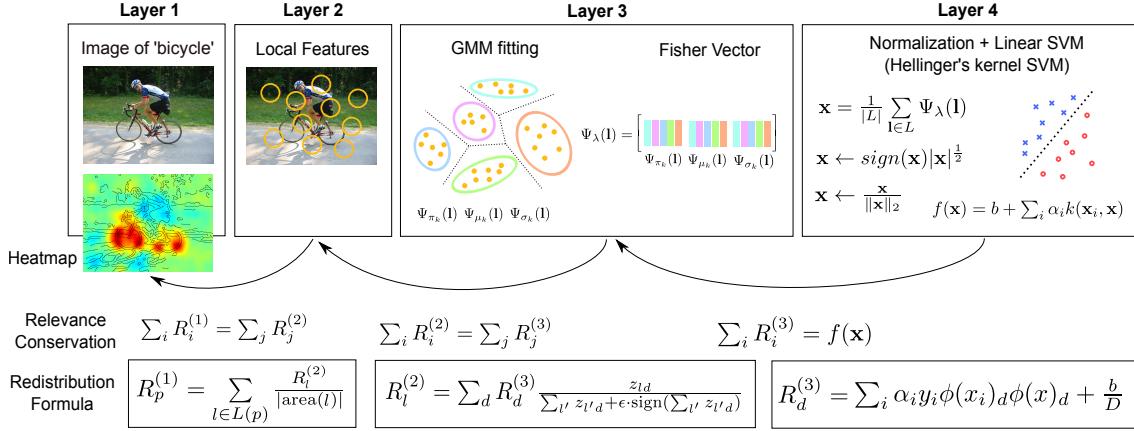
Figure 1: Computing Fisher Vector representation of an image and explaining the classification decision.

Having computed all those (as we will refer to now as) *raw* Fisher embeddings for all individual local descriptors, a single image-wise descriptor is achieved by averaging over the complete set of $\Psi_\lambda(l)$, followed by power normalization to reduce the sparsity of the descriptor and $\ell_2$-normalization to improve prediction performance [20]. The application of both final normalization steps results in a so called *improved Fisher Kernel* and is – in combination with a linear SVM [4] – equivalent to the transformation of the raw FV using the Hellinger's kernel function [20].

# 5 Explaining Classification Decisions

Most predictors, including linear SVMs over Fisher vectors, incorporate several layers of non-linear mappings, resulting in a non-linear black box with respect to the dependency of the prediction on its pixel inputs. In this section we introduce the concept of Layer-wise Relevance Propagation (LRP) [1] as a way to compute partial prediction contributions – or *relevance values R* – for intermediate and input representations based on the final classifier output. LRP acts on a single test-image similar to the work in [29] and to partial-derivative based methods such as [24]. We refer the reader to [22] for a comparison of these three explanation approaches.

## 5.1 Layer-wise Relevance Propagation

Layer-wise Relevance Propagation decomposes the mappings performed during prediction time to attribute to each component of the input its *share* with which it contributes to the classifier output, explaining its relevance to the prediction output in its given state. This unsupervised process of decomposition is in principle applicable to any kind of model, resulting in high (positive) output values $R$ identifying properties of the input speaking for the presence of the prediction target and low (or even negative) scores indicating no or negative contribution. The conservation principle inherent to LRP ensures that no amount of relevance is gained or lost in between layers of computation,

$$\sum_i R_i^{(k)} = \sum_j R_j^{(k+1)} \qquad (5)$$

where $R_i^{(k)}$ signifies the relevance value attributed to the $i$th computation unit or dimension at the $k$th computation layer of the prediction pipeline, and where the sums run over all units of the corresponding layers. In the context of an image classification problem, iterating LRP from the classifier output to the input layer results in outputs $R_p^{(1)}$ for each pixel $p$, with

$$f(x) = \sum_p R_p^{(1)} \qquad (6)$$

and $f(x)$ being equal the output layer relevance values. In [1] examples have been given for decompositions of neural network architectures and Bag of Words feature extraction pipelines satisfying the above constraints.

LRP propagates the relevance $R$ back from the output of a mapping towards its inputs. In a neural networks, a neuron maps a set of inputs $\{x_i\}$ to an output $x_j$ with monotonously increasing activation

function $g(\cdot)$

$$x_j = g(z_j) \quad \text{with} \quad z_j = \sum_i z_{ij} \ , \ z_{ij} = w_{ij} x_i$$

where the sum runs over all input neurons contributing to the activation of neuron $x_j$. The goal is to compute a relevance $R_i$ for input $x_i$ when relevances $R_j$ for outputs $x_j$ are given. [1] has introduced two possible formulas for relevance propagation

$$R_i = \sum_{j:i \to j} \frac{z_{ij}}{z_j + \epsilon \cdot \text{sign}(z_j)} R_j \qquad (7)$$

$$R_i = \sum_{j:i \to j} \left( \alpha \frac{z_{ij}^+}{z_j^+} - \beta \frac{z_{ij}^-}{z_j^-} \right) R_j, \qquad (8)$$

where $\sum_{j:i \to j}$ denotes a sum of all mappings which take $x_i$ as input. $z_{ij}^+$ denotes the positive part of the term, i.e. $\max(0, z_{ij})$, $z_j^+$ is the sum over these positive parts. $z_{ij}^-$ is defined analogously as the negative part. The same paper has introduced a method to compute relevances for Bag of words (BoW) vectors, however, it tacitly assumed that BoW mappings are dominantly non-negative. For Fisher vectors this assumption does not hold, as the features are derivatives with respect to parameters. For this reason we propose a modified approach.

## 5.2 LRP for Fisher Vector Classifiers

Our variant to use LRP for Fisher vectors starts with writing the linear SVM as a mapping of features

$$f(x) = b + \sum_i \alpha_i y_i \sum_{d=1}^{D} \phi(x_i)_d \phi(x)_d,$$

where $x$ is a raw Fisher vector, and $\phi(x)$ realizes its normalization. In consistency with the first LRP formula, we define $R^{(3)}(x)$ as

$$R_d^{(3)} = \sum_i \alpha_i y_i \phi(x_i)_d \phi(x)_d + \frac{b}{D} \qquad (9)$$

From here on we apply for the mapping of local features $l$ to Fisher vectors $x$, equation (7) instead of the approach used in [1]. We can write the d-th dimension of the Fisher vector $x_d = \sum_l m_d(l)$. This is a mapping of local features $l$ onto the Fisher vector as a set of outputs $(x_d)_{d=1}^{D}$. We apply equation (7) with $z_{ld} = m_d(l)$. $m_d(l)$ is given in the notation of Section 4 as the term from equation 13:

$$m_{(d)}(l) = \Psi_\lambda(l)_{(d)} \qquad (10)$$

Pixel-wise relevance scores $R_p^{(1)}$ are then computed by uniformly distributing for all local features $l$ the relevance scores $R_l^{(2)}$ onto the set of pixels $p$ covered by the receptive field of $l$, resulting in a *heatmap* which can be visualized. The decomposition process with explicit redistribution formulas is depicted in Figure 1.

## 5.3 Measuring Context with LRP

The distribution of positive relevance mass in a heatmap can be used for assessing the importance of context for a particular image classification task. If bounding box annotation are available (as for the Pascal VOC dataset), we can compute the *outside-inside relevance ratio* metric defined as:

$$\mu = \frac{\frac{1}{|P_{\text{out}}|} \sum_{q \in P_{\text{out}}} R_q^{(1)}}{\frac{1}{|P_{\text{in}}|} \sum_{p \in P_{\text{in}}} R_p^{(1)}} \qquad (11)$$

with $|\cdot|$ being the cardinality operator and $P_{\text{out}}$ and $P_{\text{in}}$ being the set of pixels outside and inside the bounding box, respectively. A high relevance ratio indicates that the classifier uses a lot of context to support the decision. A low relevance ratio indicates that the classifier focuses instead on the object to support its decision. Note that this measure can not be to 100% accurate in most cases, since for example the bounding box areas of slim but obliquely angled objects, for example, aeroplanes photographed during lift-off, will also cover a considerable amount of image background.

# 6 Experimental Evaluation

## 6.1 Basic Setup

All measurements are carried out on PASCAL VOC2007 [5] test data. Fisher vectors are computed using the encoding evaluation toolkit (version 1.1) from [3] with settings as in this paper. The Fisher vectors are trained on the trainval part of the same dataset. The neural network is finetuned on the trainval part of PASCAL VOC2012, starting from the BVLC reference classifier of the Caffe package [9] with a base learning rate of 0.001 using a multi-label hinge loss. As we are interested in the ability of a neural net to use context, we do not use the bounding box ground truth to extract image patches which cover parts of bounding boxes. Instead we create 4 edges and one center crop per

4

image together with mirroring, resulting in 10 training patches per image. Test scoring is done in the same fashion. This corresponds to a setting with only a few number of test windows, in which one would use larger patches during training and testing. The region-wise scores are computed for FV as described in Section 5 using equation (7) with parameter $\epsilon = 1$ and $\epsilon = 100$. For neural nets we used equation (7) with $\epsilon = 1, \epsilon = 100$ and equation (8) with $\beta = 1, \alpha = 2$. Random perturbations for Fisher vectors were achieved by randomly sampling local features from the GMM.

## 6.2 Are Fisher Explanations Meaningful ?

The first step before measuring the amount of context is to validate whether the computed scores for a pixel or a region are meaningful at all. Figure 2 depicts heatmaps computed on exemplary test images of the Pascal VOC data set considering the prediction score for a particular class. The quality of these explanations can be intuitively assessed by a human, e.g., it makes perfectly sense that the Fisher vector classifier finds that wheels are relevant for the class 'bike', rail tracks are indicative for the class 'train' and tableware is important for classifying images of class 'dining table'. These examples show that the largest part of the relevance mass does not necessarily need to lie on the object, on the contrary it may be the context which is the informative part.

In order to objectively validate that the Fisher vector heatmaps are meaningful we evaluate the decrease of the prediction score under perturbations. The idea is that a region such as an image patch is highly relevant, if modifying it results for most modifications in a sharp decline of the prediction for the whole image. Modifying a region is done by randomly perturbing the pixels with noise. The prediction score is averaged over a number of random perturbations, in order to capture the average change of the classifier.

This notion of relevant regions can be used for evaluation of region scores by sorting image regions along descending scores. Then, for each region in the sequence the average decrease of predictions is measured. The result is a graph as a function of the sequence index. Thus under this evaluation scheme, a region-wise score performs well if it assigns highest scores to regions which are most sensitive on average under perturbations and yield the sharpest decline of the prediction score. [22] introduced this setup and evaluated the methods of [29, 24, 1]

for deep neural networks tested on ImageNet [21], SUN397 [28] and MIT Places [31]. Here we show that LRP scores computed are also meaningful for Fisher vectors. Figure 7 shows this comparison against random orderings for scores computed. The LRP scores produce a more meaningful ordering than random sequences which motivates its use to define a measure for context.

## 6.3 Shallow vs. Deep Features

We investigate in the light of the LRP framework what are the differences of strategies used to classify images between (1) a shallow model operating on high-resolution images: the FV model, and (2) a deep model operating on lower-resolution images: the DNN model. We consider first the class "sheep" for which the DNN produces much better predictions than the FV model (25% superior accuracy in absolute terms according to Table 1).

Example of two images of class "sheep" and the corresponding heatmaps for the FV and DNN models are shown in Figure 4. The LRP analysis reveals that the FV and DNN models use clearly different strategies to predict the class:

The FV model bases its decision on the wool texture typical of the sheep and available at high-resolution, but ignores the exact shape of the sheep. Interestingly, relevance is also allocated to the context (here, positive relevance for the grass and negative relevance for the human face), indicating that the context is an essential component of the classifier and modulates the prediction score positively or negatively.

On the other hand, the DNN assigns a large proportion of heat to the border of the sheep, thus, showing that the shape of the sheep (e.g. its contour) is exploited in order to improve the prediction. Furthermore, for the DNN, the LRP method does not assign relevance to contextual elements such as the grass, or the human face, nor to the wool texture of the sheep, which is harder to detect due to the low resolution of images given to the CNN.

Overall, the LRP analysis indicates that the far superior predicting power of the DNN model must be attributed in largest part to the ability to model the exact shape of the sheep, making all remaining contextual or texture features less relevant. On the other hand, the less accurate FV model does benefit from the weak correlations between object class, texture and context to improve prediction quality.

|  | aeroplane | bicycle | bird | boat | bottle | bus | car |
|---|---|---|---|---|---|---|---|
| **Fisher** | 79.08% | 66.44% | 45.90% | 70.88% | 27.64% | 69.67% | 80.96% |
| **DeepNet** | 88.08% | 79.69% | 80.77% | 77.20% | 35.48% | 72.71% | 86.30% |
|  | cat | chair | cow | diningtable | dog | horse | motorbike |
| **Fisher** | 59.92% | 51.92% | 47.60% | 58.06% | 42.28% | 80.45% | 69.34% |
| **DeepNet** | 81.10% | 51.04% | 61.10% | 64.62% | 76.17% | 81.60% | 79.33% |
|  | person | pottedplant | sheep | sofa | train | tvmonitor | mAP |
| **Fisher** | 85.10% | 28.62% | 49.58% | 49.31% | 82.71% | 54.33% | 59.99% |
| **DeepNet** | 92.43% | 49.99% | 74.04% | 49.48% | 87.07% | 67.08% | 72.12% |

Table 1: Prediction performance of the trained Fisher model and deep network in average precision (AP) per class.



Figure 2: Images shown next to the heatmaps computed by application of LRP on the FV model when considering the prediction score for a particular class.
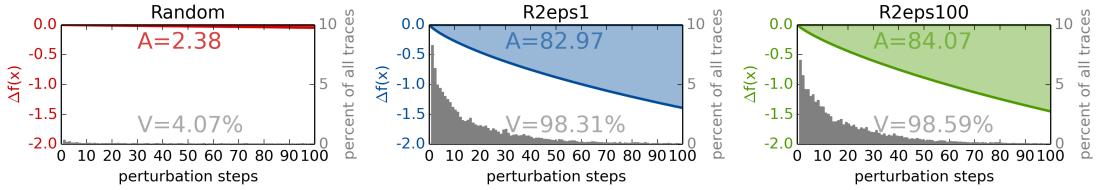
Figure 3: Heatmap Quality Measurements for Fisher Vectors. The value A measures the area above the curve between the original prediction $f(x)$ and the averaged perturbed prediction at step $i$ in the sequence of regions. $f(x) - f(x_{\text{MoRF}}^{(i)})$. V represents the fraction of all perturbation sequences for which the prediction switched sign at some step in the sequence, with the gray bar chart showing how many sample traces changed class at each point of measurement.
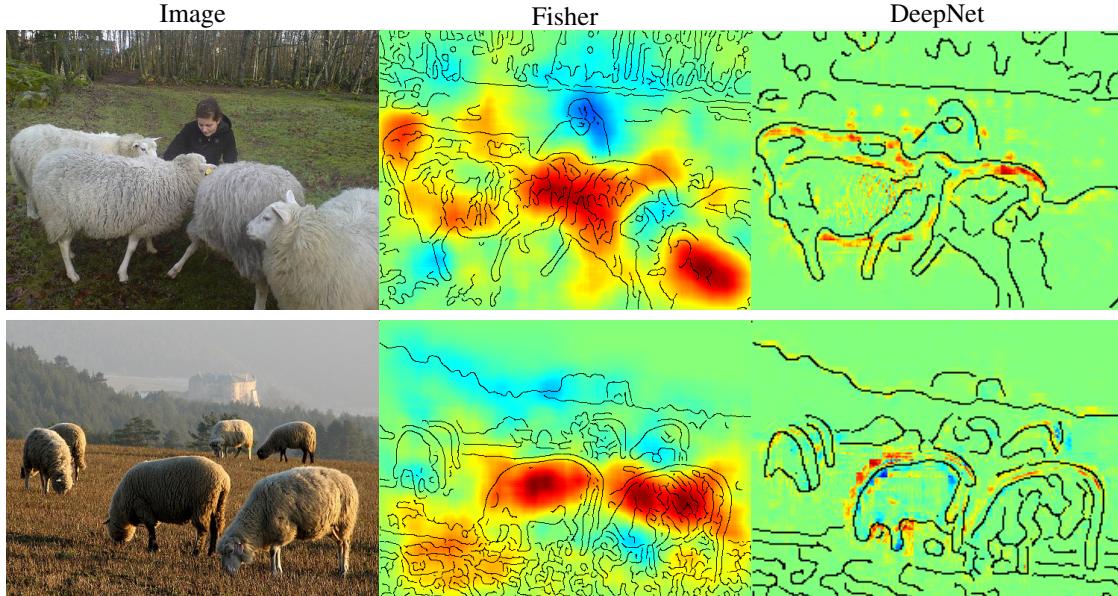


Figure 4: Images of the class "sheep", processed by the FV and DNN models and heatmapped using LRP.

## 6.4   Test Error and Model Quality

For other classes, it can be observed in Table 1 that test error of the FV model is almost on par with the one of the DNN. We investigate whether high test accuracy is predictive of the ability of the model to extract meaningful features for a given class, or whether the decision is based mostly on undesirable contextual or artefactual features.

**Contextual features**   As an illustrative example, we consider the class "boat", where the performance of the DNN superior by less than 7% in absolute terms to the FV model. (Note that for other classes such as "sheep" or "bird", the DNN performance is superior by 25% or more.) It is tempting to conclude that, for the class "boat", both models should have learned a set of features of similarly high quality.

LRP analysis gives a different answer: Figure 5 (left) shows the heatmaps produced by the FV and DNN models on two archetypical images of the class "boat". For the DNN, LRP assigns most of the relevance to pixels corresponding to the actual boat. On the other hand, for the FV model, LRP assigns most relevance to the water below the boat (i.e. the FV model does not recognize the object itself, but its context). The heat distribution of average heatmaps (computed over all landscape-format images of the class "boat") corroborates what was observed for two selected images, in particular, a focus of the FV model on the bottom part of the image where water usually is, and a focus of the DNN model on the middle part of the image where the boat typically is.

We can conclude from the LRP analysis, that while both classifiers have a roughly similar level of
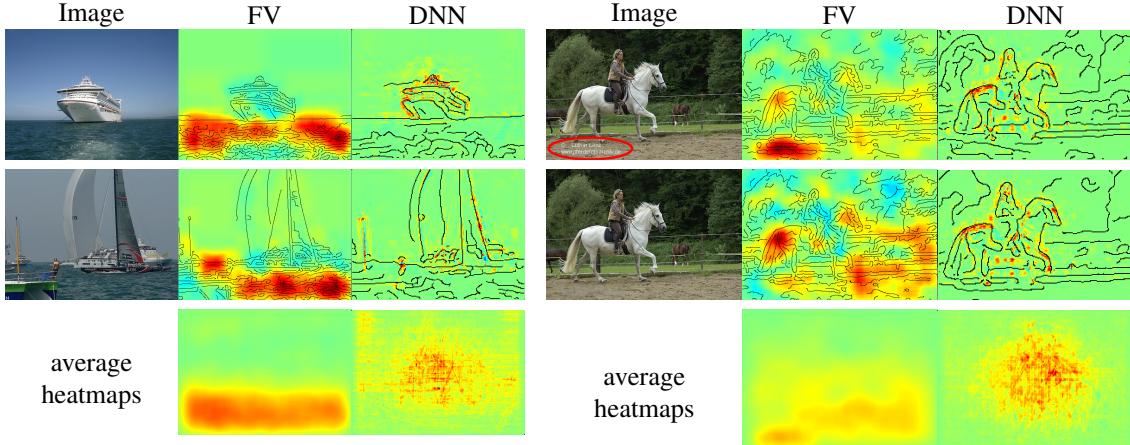
| Image | FV | DNN | Image | FV | DNN |
|---|---|---|---|---|---|



Figure 5: Top: Images of the classes "boat" and "horse", processed by the FV and DNN models and heatmapped using LRP. Bottom: Average heatmap scores over a random sample (of size between 47 and 177) of the distribution for each class and model. On the second image of class "horse", the copyright tag (marked by the red ellipse) has been removed.

accuracy on the test images with class "boat", FV's performance is likely to decrease drastically if one were to consider boats located outside the water as test images. On the other hand, performance of the DNN would be less affected. Therefore, test error is a superficial predictor of model quality in this case.

**Artefactual features**  A second example where high accuracy does not necessarily translates into high quality features is for the class "horse". This class is predicted with similar accuracy by the FV and DNN models (approximately 1% difference in accuracy).

Figure 5 (right) shows a LRP heatmaps for the FV and DNN model on an image of horse. While the DNN assigns relevance on the actual object "horse", the FV assigns almost all relevance in the bottom-left corner the image, where careful inspection of the image reveals the presence of a copyright tag. Thus, the decision of the FV model is in large part based on the presence of the copyright tag, which is discriminative of the class horse. Removing the copyright tag completely changes the FV heatmap, but does not change significantly the DNN heatmap.

If the copyright tag is removed, the DNN is still able to predict the image because the pixels that support its decision are not affected. On the other hand, FV model prediction quality will be considerably reduced. The systematic focus of the FV model on the copyright tag is confirmed in the average heatmap, where the bottom-left corner is assigned large amount of heat. Therefore, for this class

again, test error does not predict well model quality.

## 6.5  Quantitative analysis of context use

While we have so far provided a qualitative interpretation of FV heatmaps for examples and classes of interest, we can more systematically measure whether the model uses context or the actual object, by measuring for each classes and models the outside-inside relevance ratio $\mu$ computed by equation 11.

Results are shown in Figure 6. Generally, the FV model uses more context than the DNN, as evidenced by a higher relevance ratio. However, there are significant differences between classes:

Classes where the use of context by the FV model is particularly high are "boat" and "airplane", the first of which we have studied qualitatively in the previous section. For these two respective classes, the water and the sky are important contextual elements that support the decision of the Fisher model, due to their strong correlation. Another group of classes with high context of the Fisher model are "chair", "diningtable", "pottedplant" and "sofa" which share a semantic of indoor room sceneries.

For other classes such as "bicycle", "car", "motorbike", or "sheep", the Fisher model does not use much context. For the first three classes, the urban environment surrounding these classes is not predictive of the object being detected (i.e. it could not discriminate between these three classes based on the context only). For the last class, as it has been discussed in Section 6.3, the wool texture of the sheep
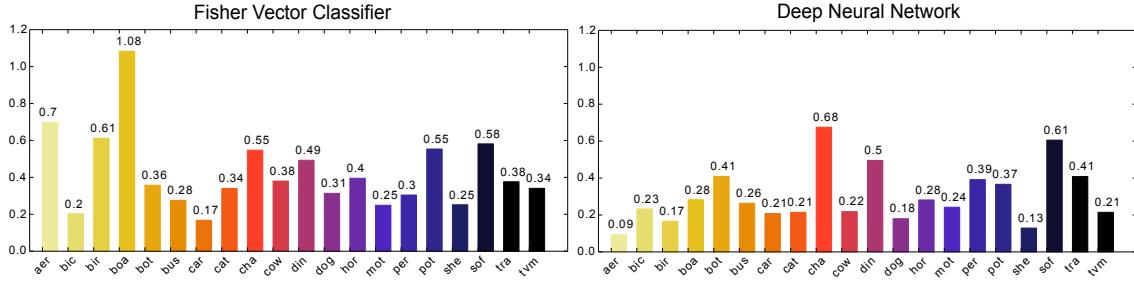
Figure 6: Outside-inside relevance ratio as computed by equation 11 for the 20 classes of the Pascal VOC 2007 dataset. Left: ratios for the FV model. Right: ratios for the DNN model.

(which lies inside the sheep bounding box) is a reasonable predictor for the class "sheep", although the actual object sheep (i.e. defined by its shape or contour) is not being used.

As for deep neural networks, classes with least context usage are "aeroplane", "bird", "sheep", "dog", "car", "cat" and "tvmonitor". Each of those is associated with a significantly better score achieved by the DNN.

## 7  Conclusion

In this paper, we have analyzed what make Fisher vector models (FV) and deep neural networks (DNN) decide for a particular class. To achieve this, we have employed a heatmapping technique that determines what pixels in the image are used by a classifier to support its decision. The technique called layer-wise relevance propagation and originally developed for neural networks [1] was extended to Fisher vector models, and validated using the method by [22].

Our novel comparative analysis of FV and DNN classifiers corroborates empirically previous intuition relating the architecture of the classifier to the features it is able to extract. In particular, our analysis shows that the FV model compensates its lack of depth by the use of contextual information—potentially artefacts—that are weakly correlated to the object class. We thus demonstrate that the generalization capability of Fisher vector models can be overstated if test images also include similar context.

On the other hand, DNNs base their decision on the actual object to detect and ignores its context. This focus on object detection has to be attributed to the higher overall predictive accuracy of the model, that removes the need for contextual information—even if the latter is discriminative. The focus on detection must also be attributed to the deep multi-

task properties of the DNN that favors composition of natural image features over lower-level features such as copyright text.

These results argue in favor of incorporating heatmapping techniques into the data collection and model selection processes. The interpretable visual feedback that heatmaps provide can be used in particular to verify that the considered classifier bases its decision on the right set of features, and in the contrary case, select another model, or extend the dataset in a way that artefactual features can no longer support the classification decision.

## References

[1] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS ONE*, 10(7):e0130140, 2015.

[2] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K.-R. Müller. How to explain individual classification decisions. *JMLR*, 11:1803–1831, 2010.

[3] K. Chatfield, V. S. Lempitsky, A. Vedaldi, and A. Zisserman. The devil is in the details: an evaluation of recent feature encoding methods. In *BMVC*, page 8, 2011.

[4] C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.

[5] M. Everingham, S. M. A. Eslami, L. V. Gool, C. K. I. Williams, J. M. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *IJCV*, 111(1):98–136, 2015.

[6] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results.

[7] R. B. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. *CoRR*, abs/1311.2524, 2013.

[8] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, 1997.

[9] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.

[10] A. Karpathy, A. Joulin, and F. Li. Deep fragment embeddings for bidirectional image sentence mapping. In *Adv. in NIPS*, pages 1889–1897, 2014.

[11] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and F. Li. Large-scale video classification with convolutional neural networks. In *CVPR*, pages 1725–1732, 2014.

[12] J. Koutník, G. Cuccu, J. Schmidhuber, and F. J. Gomez. Evolving large-scale neural networks for vision-based reinforcement learning. In *GECCO*, pages 1061–1068, 2013.

[13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Adv. in NIPS*, pages 1106–1114, 2012.

[14] Y. LeCun and C. Cortes. The MNIST database of handwritten digits. http://yann.lecun.com/exdb/mnist/, 1998.

[15] K. Lenc and A. Vedaldi. Understanding image representations by measuring their equivariance and equivalence. *CoRR*, abs/1411.5908, 2014.

[16] L. Liu and L. Wang. What has my classifier learned? visualizing the classification rules of bag-of-feature model by support region detection. In *CVPR*, pages 3586–3593. IEEE, 2012.

[17] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. *CoRR*, abs/1412.0035, 2014.

[18] D. Novotný, D. Larlus, F. Perronnin, and A. Vedaldi. Understanding the fisher vector: a multimodal part model. *CoRR*, abs/1504.04763, 2015.

[19] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *CVPR*, pages 1717–1724, 2014.

[20] F. Perronnin, J. Sánchez, and T. Mensink. Improving the fisher kernel for large-scale image classification. In *Computer Vision–ECCV 2010*, pages 143–156. Springer, 2010.

[21] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, pages 1–42, April 2015.

[22] W. Samek, A. Binder, G. Montavon, S. Bach, and K. Müller. Evaluating the visualization of what a deep neural network has learned. *CoRR*, abs/1509.06321, 2015.

[23] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek. Image classification with the fisher vector: Theory and practice. *IJCV*, 105(3):222–245, 2013.

[24] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *ICLR Workshop*, 2014.

[25] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Av. in NIPS*, pages 3104–3112, 2014.

[26] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014.

[27] J. Uijlings, A. Smeulders, and R. Scha. The visual extent of an object. *IJCV*, 2012.

[28] J. Xiao, J. Hays, K. Ehinger, A. Oliva, A. Torralba, et al. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, pages 3485–3492. IEEE, 2010.

[29] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, pages 818–833, 2014.

[30] N. Zhang, J. Donahue, R. B. Girshick, and T. Darrell. Part-based r-cnns for fine-grained category detection. *CoRR*, abs/1407.3867, 2014.

[31] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *Adv. in NIPS*, pages 487–495, 2014.

# A    Details on the Computation of Fisher Vector Embeddings

The representation of an image data set as FVs starts with the computation of a Gaussian Mixture Model (GMM) as a soft vocabulary of visual prototypes. The $K$ components $\lambda = \{(\pi_k, \mu_k, \Sigma_k)\}_{k=1..K}$ of the GMM are fitted onto a set of local descriptors $l$ extracted from the training images, where $\mu_k$ is the mean vector of the $k$th mixture component and $\Sigma_k$ its covariance matrix which is assumed to be diagonal, e.g. $diag(\Sigma_k) = \sigma_k$ and all off-diagonal entries are 0. The parameter $\pi_k$ is the mixture weight of component $k$, with $\sum_k \pi_k = 1$ and $\forall k : \pi_k \geq 0$. A *full* FV descriptor of an image measures the average 0th (soft mapping weight), 1st (deviation from mean) and 2nd (variance) moment of all local descriptors sampled from the image area [23].

$$\Psi_{\pi_k}(l) = \frac{1}{\sqrt{\pi_k}} \left(\gamma_k(l) - \pi_k\right)$$

$$\Psi_{\mu_k}(l) = \frac{1}{\sqrt{\pi_k}} \gamma_k(l) \left(\frac{l - \mu_k}{\sigma_k}\right)$$

$$\Psi_{\sigma_k}(l) = \frac{1}{\sqrt{\pi_k}} \gamma_k(l) \frac{1}{\sqrt{2}} \left(\frac{(l - \mu_k)^2}{\sigma_k^2} - 1\right)$$

$$(12)$$

with $\Psi_{\pi_k}(l) \in \mathbb{R}$, $\{\, l,\ \Psi_{\mu_k}(l),\ \Psi_{\sigma_k}(l)\,\} \in \mathbb{R}^D$ and $\gamma_k(l)$ returning the soft assignment of $l$ to the $k$th mixture component [23]. The FV embedding $\Psi_\lambda(l)$ of a local descriptor $l$ is achieved by concatenating the mapping outputs for all $K$ components into a $(1 + 2D)K$ dimensional vector, such that

$$\Psi_\lambda(l) = \left[\underbrace{\Psi_{\pi_1}(l),\ \ldots,\ \Psi_{\pi_K}(l)}_{K\ \text{scalar entries}},\quad \underbrace{\Psi_{\mu_1}(l),\ \ldots,\ \Psi_{\mu_K}(l)}_{K\ \text{concatenated}\ D\text{-dimensional vectors}},\quad \underbrace{\Psi_{\sigma_1}(l),\ \ldots,\ \Psi_{\sigma_K}(l)}_{K\ \text{concatenated}\ D\text{-dimensional vectors}}\right] \quad (13)$$

The final FV of an image is obtained by averaging over all $\Psi_\lambda(l)$ resultung in the *raw* FV representation[1], followed by power normalization to reduce the sparsity of the descriptor and and $\ell_2$-normalization to improve prediction performance [20]:

$$x = \frac{1}{|L|} \sum_{l \in L} \Psi_\lambda(l) \qquad\qquad \text{mapping aggregation}$$

$$x \leftarrow sign(x)|x|^{\frac{1}{2}} \qquad\qquad \text{power normalization}$$

$$x \leftarrow \frac{x}{\|x\|_2} \qquad\qquad \ell_2\text{-normalization}$$

# B    On the Equality of the Normalization Steps and the Hellinger's Kernel

The work of [23] references the equality of the application of the Hellinger's (Bhattacharyya) kernel function to a raw FV descriptor and the improved FV (power- and $\ell_2$-normalized FV) with a linear kernel function on top. Here, we explicitly show this equality. Assume a raw FV $x$, the component-wise absolute of an input variable $|\cdot|$, and the $\ell_p$-norm of a vector $\|\cdot\|_p$. The improved FV defines itself as the application of power-normalization followed by $\ell_2$-normalization on top of a raw FV, e.g. interpreted as a mapping

---

[1] *raw* FV in contrast to the *improved* FV, as used in the main document to distinguish between the both.

function $\Phi(\cdot)$ we receive

$$
\Phi(x) \;=\; \frac{sign(x)|x|^{\frac{1}{2}}}{\|sign(x)|x|^{\frac{1}{2}}\|_2} \;=\; \frac{sign(x)\sqrt{|x|}}{\sqrt{\sum_d \underbrace{sign(x_{(d)})^2}_{=1} \underbrace{|x_{(d)}|^{\frac{1}{2}\cdot 2}}_{=1}}} \tag{14}
$$

$$
=\; sign(x)\sqrt{\frac{|x|}{\sqrt{\sum_d |x_{(d)}|}^2}} \;=\; sign(x)\sqrt{\frac{|x|}{\|x\|_1}}
$$

and with that

$$
k(x, \boldsymbol{y}) = \sum_d \Phi(x)_d \Phi(\boldsymbol{y})_d = \sum_d sign(x_{(d)}y_{(d)}) \underbrace{\sqrt{\frac{|x_{(d)}|}{\|x\|_1} \cdot \frac{|y_{(d)}|}{\|\boldsymbol{y}\|_1}}}_{\text{Hellinger's kernel}} \tag{15}
$$

which enables us to compute component-wise relevance scores for each FV dimension $d$ by attributing the power normalization to the kernel function as

$$
R_d^{(3)} = \sum_i \alpha_i \Phi(x_i)_d \Phi(x)_d + \frac{b}{D} \tag{16}
$$

In the case of an improved FV mapping (or the Hellinger's (Bhattacharyya) kernel function), the support vectors need not to be known explicitly, since $w = \sum_i \alpha_i \Phi(x_i)$, further simplifying Equation to

$$
R_d^{(3)} = w_d \Phi(x)_d + \frac{b}{D} \tag{17}
$$

# C   Details on Relevance Decomposition for Fisher Vectors

As already stated within the main document, the relevance decomposition for FV embeddings follows the definitions and constraints of [1]. As already stated in Section B, relevance values $R_d^{(3)}$ for each dimension $d$ of the FV representation of an image can be easily computed as

$$
R_d^{(3)} = \sum_i \alpha_i \Phi(x_i)_d \Phi(x)_d + \frac{b}{D} \tag{18}
$$

The next step towards local explanations in pixel space is the computation of relevalce scores $R_l^{(2)}$ for each local descroptor $l$ contributing to the computation of the FV used for prediction. The work of [1] introduces the notion of a mapping function $m_{(d)}(l)$ relating $l$ to dimension $d$ in output space for Bag of Feature (BoF) models assuming positive mapping outputs exclusively. Reformulating $\Psi_\lambda(l)$ in terms of $m_{(d)}(l)$ to fit it into the LRP framework facilitates the computation of local feature relevance scores proportionally to their forward mapping contributions

$$
m_{(d)}(l) = \begin{cases} \frac{1}{\sqrt{\pi_k}}(\gamma_k(l) - \pi_k) & ; d = k, \ k \in [1, K] \\ \frac{1}{\sqrt{\pi_k}}\gamma_k(l)\left(\frac{l_{(r)} - \mu_{k,(r)}}{\sigma_{k,(r)}}\right) & ; d = K + D(k-1) + r, k \in [1, K], r \in [1, D] \\ \frac{1}{\sqrt{\pi_k}}\gamma_k(l)\frac{1}{\sqrt{2}}\left(\frac{(l_{(r)} - \mu_{k,(r)})^2}{\sigma_{k,(r)}^2} - 1\right) & ; d = (1 + D)K + D(k-1) + r, k \in [1, K], r \in [1, D] \end{cases} \tag{19}
$$

With that the proposed decomposition formula for BoF mappings can directly be applied with only a minor adaption of $Z(x)$ from Equation (26) in [1] to consider that $m_d(l)$ may output mapping weights of both positive and negative signs for FVs :

$$
R_l^{(2)} = \sum_{d \notin Z(x)} R_d^{(3)} \frac{m_{(d)}(l)}{\sum_{l' \in L} m_{(d)}(l')} + \xi \tag{20}
$$

where

$$Z(x) = \left\{ d \mid \forall l : m_{(d)}(l) = 0 \right\} \text{ and } \xi = \sum_{d \in Z(x)} R_d^{(3)} \frac{1}{|L|}$$

Pixel-wise relevance scores $R_p^{(1)}$ are then computed by uniformly distributing for all local features $l$ the relevance scores $R_l^{(2)}$ onto the set of pixels $p$ covered by the receptive field of $l$

$$L(p) = \{ l \mid p \in \text{area}(l) \}$$

$$R_p^{(1)} = \sum_{l \in L(p)} \frac{R_l^{(2)}}{|\text{area}(l)|} \tag{21}$$

resulting in a *heatmap* which can be visualized.

Above approach to compute values $R_l^{(2)}$ can in general – when $m_{(d)}(l)$ outputs values of both signs – be expected to be numerically problematic when $\sum_{l'} m_{(d)}(l')$ becomes very small due to individual mappings $m_{(d)}(l')$ cancelling each other out. Divisions close to zero would then pronounce otherwise insignificant local descriptor weights for relevance distribution. This problem is known in the context of neural network relevance decompositions and has been discussed in [1], which we will follow for extending Equation 20 to better satisfy the task at hand. Figure 7 will then compare how well the computed heatmap resulting from each decomposition method represents the classifier's perception according to the measurement procedure described in Section D.
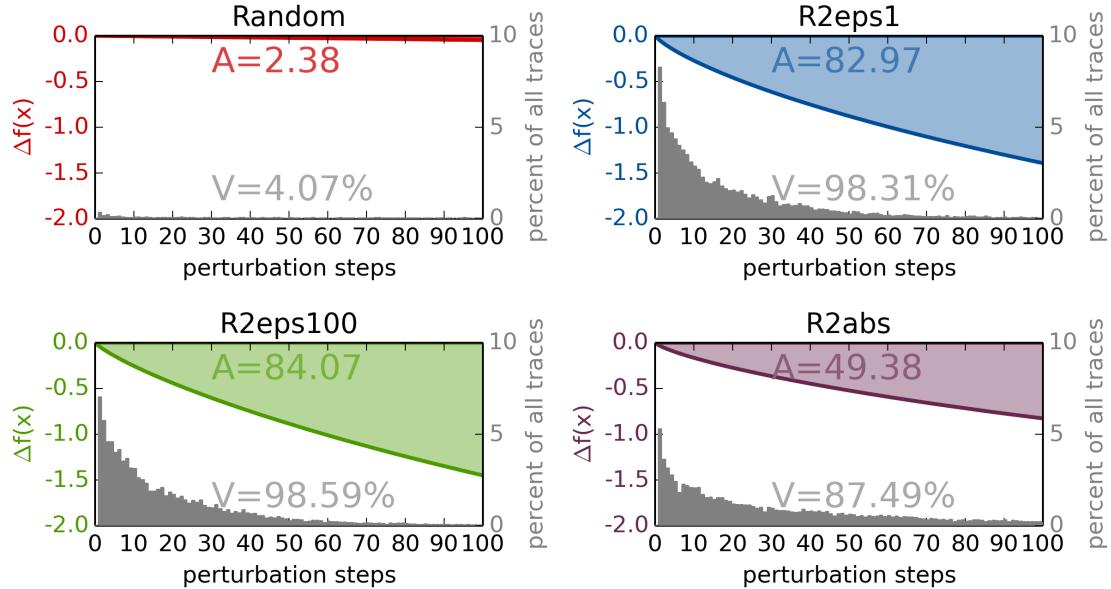


Figure 7: Heatmap Quality Measurements for Fisher Vectors. The value A measures the area above the curve between the original prediction $f(x)$ and the averaged perturbed prediction at step $i$ in the sequence of regions. $f(x) - f(x_{\text{MoRF}}^{(i)})$. V represents the fraction of all perturbation sequences for which the prediction switched sign at some step in the sequence, with the gray bar chart showing how many sample traces changed class at each point of measurement. Despite violating the relevance conservation principle we have found the $\epsilon$-stabilized decomposition to produce the best heatmaps for the FV classifer. Therefore heatmaps computed with Equation 22 and $\epsilon = 100$ for all further evaluations.

## C.1   $\epsilon$-Stabilized Decomposition

The first and probably most straight forward adaption of the decomposition formula introduced in [1] is the introduction of a numerical stabilizer $\epsilon > 0$ to prevent (near) zero divisions

$$R_l^{(2)} = \sum_{d \notin Z(x)} R_d^{(3)} \frac{m_{(d)}(l)}{\sum_{l' \in L} m_{(d)}(l') + \gamma} + \xi \tag{22}$$

with $\gamma = \epsilon \cdot sgn\left(\sum_{l' \in L} m_{(d)}(l')\right)$, $sgn(y) = 1$ if $y \geq 0$ and $-1$ else. An obvious drawback of this approach is, however, that it violates the conservation constraints set by LRP, as varying amounts of relevance can be absorbed or generated by $\epsilon$.

## C.2   Absolute Value Decomposition

Our second alternative to resolve the issue of mapping weight cancellation is to only consider the absolute of a descriptor's mapping contribution, effectively removing all cases of division by zero not already covered by Equation 20. As a further benefit, no additional parameters need to be defined and no relevance is lost as in Eq. 22

$$R_l^{(2)} = \sum_{d \notin Z(x)} R_d^{(3)} \frac{|m_{(d)}(l)|}{\sum_{l' \in L} |m_{(d)}(l')|} + \xi \tag{23}$$

The drawback of this method is the loss of information encoded in the mapping signs.

# D   Measuring the Quality of a Heatmap

With *pixel flipping* a procedure has been introduced in [1], which demonstrated that the heatmaps computed for the MNIST [14] data set are able to successfully identify the properties of the input image determining the decision of the classifier. The idea behind this data deterioration approach guided by relevance scores has then been extended in [22] to an evaluation procedure to compare alternatives for LRP decomposition and image perturbation approaches applied to a deep neural network classifier.

Assessing the quality of a heatmap using image deterioration guided by the heatmap itself follows the semantic behind LRP: Let us assume a classifier which predicts a learned target if $f(x) > 0$. The application of LRP to $f(\cdot)$ with an input point $x$ should then attribute (the highest) positive values $R_i^{(k)}$ to components of $x$ which are (most) important to the positive decision of $f(\cdot)$. By intuition, removing that information should cause the predictor output to become less positive. Considering the values $R_i^{(k)}$ as a *ranking* of importance of components of the input , e.g. the most important components receive the largest proportions of relevance, then this ranking is optimal if the output of $f(\cdot)$ decreases faster than any other ranking with ongoing removal of input information with a removal order $\vec{o}$ determined by ordering the values $R_i^{(k)}$ descendingly. Since different decomposition methods (e.g. randomly picked component orders) can produce different orderings of importance, an algorithm of relevance-guided data deterioration can be used to compare how well each decomposition variant explains the reasoning of the classifier. We formally define this deterioration algorithm as

$$x_{\text{MoRF}}^{(0)} = x$$
$$\forall \, 1 \leq i \leq I \, : \, x_{\text{MoRF}}^{(i)} = g\left(x_{\text{MoRF}}^{(i-1)}, \vec{o}_i\right) \tag{24}$$

with $g(\cdot)$ being a perturbation function to remove or exchange data, $x_{\text{MoRF}}^{(i)}$ representing a data point after $i$ steps of removal of the **mo**st **r**elevant information **f**irst, determined by an ordering $\vec{o}$, have been performed. We compute a measure of quality of a heatmap $A$ by comparing the prediction $f(x)$ to the prediction on the altered input $x_{\text{MoRF}}^{(i)}$ and then computing the area between those two curves by integrating over the points of measurement:

$$A = \frac{1}{I} \sum_{i=1}^{I} \left(f(x) - f(x_{\text{MoRF}}^{(i)})\right) \tag{25}$$

14

The value $A$ increases with the reaction of the classifier to a change in its input and therefore higher values for $A$ encode a better representation of the classifier decision in terms of back propagated relevance.

## D.1 Local Feature Replacement

Previously used variants of pixel flipping have – hence the name – operated on pixel level exclusively to evaluate heatmaps for neural network type classifiers. Those classifiers – namely deep neural networks – directly received the pixel values of an image as inputs and the application of LRP produce pixel-accurate heatmaps. While [1] exchanges almost binary pixel values of the $28 \times 28$ pixel sized MNIST pixels by inverting the state of each pixel, [22] studies different image perturbation strategies on $227 \times 227$ pixel large color images showing photographic scenes and structures. Here, questions such as the number of pixels to exchange, their spatial grouping and how to best replace those pixels have been raised and the overall complexity of the problem of finding the *right* perturbation strategy is discussed.

When a feature extraction pipeline (in the classical sense: dense local feature extraction $\rightarrow$ mapping $\rightarrow$ pooling) is part of the predictor operating on images we face a set of distinct problems. For example, exchanging one pixel causes a recalculation and mapping of all local descriptors covering that pixel in order to measure the effect to the predictor output. Assuming a meaningful heatmap, pixel positions with leading relevance scores are in general not sparsely scattered all over the image area but rather grouped in the same image area over the extend of a patch of neighbouring pixels. This is even more so true due to the computation of pixel relevance scores $R_p^{(1)}$. The same local descriptors are very likely to be recalculated over and over again, even though the change in a single pixel will not have much of an effect. This causes a considerable computational cost to assess a heatmap, especially when local descriptors are sampled densely and at multiple scales. An obvious solution to this problem is to exchange a group of pixels at a time, yet this raises the questions of how many pixels to exchange, due to which grouping and with what replacement strategy? Removing single pixels is out of the question, since this would effectively damage the integrity of the image itself. Also the type of local descriptor needs to be considered, e.g. do SIFT features encode shape information well. By blurring an area important structural information might be removed, and a bias towards another class (or even the image's true class) might be introduced, by the blurred area resembling e.g. sky, water or cloth, interfering with the evaluation for classes *aeroplane*, *boat* and *person*, respectively. Vice versa, setting selected pixels to random color values might create high contrast gradients in the image, affecting the evaluation of certain object classes identifying themselves via sharp and well pronounced edges.

We aim to avoid all those problems by not exchanging pixel values, but local descriptors instead. Firstly, we assume the descriptors forming the orderless descriptor set $L$ serving as an intermediate image representation to be the result independent events which can also be exchanged one at a time without affecting each other's meaning. This is not the case when individual pixels are exchanged. Secondly, in contrast to pixel scores $R_p^{(1)}$ the local feature scores $R_l^{(2)}$ are also the direct result of the relevance decompositions of the mapping function $m_d(l)$ of which we wish to compare the alternative options. Se can use the GMM $\lambda$ representing the distribution of local descriptors as a generative model to draw (already dimensionality-reduced) replacements for the features to be exchanged, resulting in believable data points $x_{\text{MoRF}}^{(i)}$ close to the data manifold. What remains to do is to create an ordering of local descriptors to replace and to update the FV $x$. Algorithm 1 outlines this perturbation and evaluation strategy. Note, that the transformation $\Phi(\cdot)$ is applied to $x_{\text{MoRF}}^{(i)}$ before feeding it as an input to $f(\cdot)$ to compute $A$ in Equation 24. To measure the

---

**Algorithm 1** Local Feature Resampling

---

1: **Input:** local features $L$, GMM $\lambda$, number of replacements $I$, replacement order $\vec{o}$
2: $x_{\text{MoRF}}^{(0)} \leftarrow \frac{1}{|L|} \sum_{l \in L} \Psi_\lambda(l)$
3: **for** $i$ from 1 to $I$ **do**
4: $\quad l' \leftarrow$ feature to replace $L(\vec{o_i})$
5: $\quad l_\lambda \leftarrow \lambda.\text{generate}()$
6: $\quad x_{\text{MoRF}}^{(i)} \leftarrow x_{\text{MoRF}}^{(i-1)} + \frac{1}{|L|} \Psi_\lambda(l_\lambda) - \frac{1}{|L|} \Psi_\lambda(l')$
7: **end for**

---

impact of Algorithm 1 when applied to an image, we replace the first $10,000$ (wrt to the ordering $\vec{o}$) local descriptors extracted from the image in batches of 100, resulting in 100 points of measurement. As a rough estimate, from about $15,000$ up to $100,000$ features are computed for each image of the PASCAL VOC 2007 test set using the dense sampling approach of the encoding evaluation toolbox [3], with $69,000$ local features being sampled per image on average.

To reduce random effects influencing the measurement process, we repeat the experiment five times, in total recording 28870 image perturbation traces on all true positive predictions [2] after optimizing the prediction threshold wrt to the EER measure.

# E   Details on the Neural Network Retraining

The starting point was the BVLC reference caffe net as provided with the caffe package [9]. Training mode was multi-label training instead of the usual competitive multi-class training because for PASCAL VOC multiple classes can be present in one image. The training criterion was the sum of hingelosses over all 20 classes in the Pascal VOC data set. Note that this required to use a customized image data layer and a customized hinge loss layer.

One general problem for training and testing is the question how to score and image and what data to use for training. A second problem is how to generate patches matching the quadratic receptive field size from non-quadratic images. One general approach to generate non-quadratric images is to ignore the aspect ratio and to use warping in order to transform a non-quadratic patch into a quadratic one such as in [7].

In order to have maximal comparability to Fisher vectors which do not use warping and process an image as a whole we decided for a simpler setup which is close to the setup used by Fisher vectors and which preserved the aspect ratio of patches used during training and testing, irrespective of the fact that other setups may have resulted in somewhat higher performance of the neural network.

**Training data:**

As we were interested in training a setup such that the neural network is able to use context, we refrained from training the network with image patches around the scale of a bounding box and smaller. We decided not to use the information about object bounding boxes for generating training data because for Fisher vectors this information was not used for generating training data. Instead each image was rescaled such that the largest side had 256 pixels. The smaller side was padded at its boundaries by the nearest pixel. From this modified image 4 edge and one center crop was taken. After mirroring the image, this was repeated. This resulted in 10 images per training image. This is a compromise to ensure a sufficiently large sample size for retraining, as it is known that neural networks excel typically at higher training sample sizes.

**Testing data:**

The heatmap was computed using one center crop only.

As for measurement of mean average precision, results depend on how to score one image at test time. Note that it is common to compute an average score over many crops of the image.

Resizing the largest side of the image to 256 pixels and using the $227 \times 227$ center crop only for each image resulted in the 72.12 mAP reported in the main paper. Note that this setup corresponded to the setup used for computing the heatmaps, so this was used for the sake of comparability.

Using a different test strategy, namely resizing the smallest side of the image to 256 pixels, then computing an average over a sliding window with stride of 20 pixels, resulted in an increase of ranking performance to 75.9 mAP. This strategy used merely 12 - 35 test patches per image. Using approaches with several hundred test windows, such as 500 windows in [19] would probably have resulted in much higher mAP scores, however we did not consider a higher score relevant for the main message of the paper focused on context

---

[2] 5774 TP predictions in total over all 20 classes with 5 repetitions for each case. Single images may show multiple and correctly detected object classes at once.

usage. In view of the considerably increased computation time for such approaches we refrained from them.