

Understanding the Fisher Vector: a multimodal part model

David Novotný¹
novotda8@fel.cvut.cz

Diane Larlus²
diane.larlus@xrce.xerox.com

Florent Perronnin²
florent.perronnin@xerox.com

Andrea Vedaldi³
<http://www.vlfeat.org/~vedaldi>

¹Center for Machine Perception
Czech Technical University

²XRCE
Xerox Research Centre Europe
Meylan, France

³Department of Engineering Science
Oxford University

Abstract

Fisher Vectors and related orderless visual statistics have demonstrated excellent performance in object detection, sometimes superior to established approaches such as the Deformable Part Models. However, it remains unclear how these models can capture complex appearance variations using visual codebooks of limited sizes and coarse geometric information. In this work, we propose to interpret Fisher-Vector-based object detectors as part-based models. Through the use of several visualizations and experiments, we show that this is a useful insight to explain the good performance of the model. Furthermore, we reveal for the first time several interesting properties of the FV, including its ability to work well using only a small subset of input patches and visual words. Finally, we discuss the relation of the FV and DPM detectors, pointing out differences and commonalities between them.

1 Introduction

Object detection is a key task in computer vision and, as such, the community has dedicated to this problem a tremendous amount of effort. In the past several years, a predominant line of work in this area has been the use of sliding window classifiers computed on top of HOG [5] or similar features. The original HOG classifier of Dalal & Triggs has since been extended in several ways. The most representative of such extensions is the *Deformable Part Model* (DPM) of Felzenszwalb *et al.* [9] that constitutes today the *de-facto* standard solution for a generic object detector. The DPM improves on the basic HOG-based detector in two key ways: by the introduction of *deformable parts*, allowing templates to deform, and by the introduction of *multiple aspects*, allowing to capture very different object appearances.

A property of HOG-based models including the DPM, that may explain their popularity nearly as much as their excellent performance, is the fact that the learned models are often easily interpretable. Already [5, 9] used the structure of the HOG features to generate model visualizations that convincingly and

intuitively capture well-defined object parts. More recently, techniques such as HOGgles [19] have been proposed to take these visualizations to the next level, revealing many interesting features of these models and of their failure modalities.

HOG-like features are not the only popular image representations in computer vision. Before the introduction of HOG, image-based modeling was flourishing using orderless statistics, starting with the Bag-of-Visual-Words (BoVW) representation [4, 16]. Similar models are still a popular choice for image classification and retrieval tasks. In fact, at the same time as DPM became the most popular model for object detection, orderless statistics were found to work as well or even better than the DPMS [18] in international benchmarks such as the PASCAL VOC challenge [6]. Similar to HOG and DPM, orderless models have ever since been significantly improved; the best current representative methods include VLAD [10] and, in particular, the *Fisher Vector* (FV) [14] which, compared to BoVW, capture significantly richer statistics of the visual word occurrences. Recently, [2, 3] achieved state-of-the-art object detection performance by using FVs.

An open challenge in orderless models is understanding the nature of the visual information that they capture. Part of this challenge is that, differently from HOG, models such as BoVW, VLAD, and FV are difficult to visualize, so that it is unclear what aspects of an image or object class they model and how. The reason is that, while HOG pools local information at well-defined spatial locations in an image, orderless models scramble this information into a bag, making it difficult to reconstruct the object being recognized. Furthermore, while in DPM it is easy to define and visualize a notion of a semantic *part*, the statistical analogous in the case of BoVW, VLAD, or FV is much less clear. The goal of this paper is to shed light on these issues by understanding, interpreting and visualizing these orderless models. Our focus is the FV representation in the context of detection, but our conclusions should extend to related models such as BoVW and VLAD.

Our contributions are threefold. First, we show that the FV detector can be formulated as a part-based model where what we define as parts have a similar role to the one of parts in DPMS. Second, based on this new formulation of the FV detector, we discuss the similarities and the differences between the FV detector and the DPM, and explain how orderless models manage to capture varying object appearances using a small visual codebook. In particular, while DPM uses unimodal, movable, and well localized parts, the FV detector uses a fixed set of parts, but is capable of capturing complex multi-modal appearances of each part. Finally we demonstrate the sparsity of FV detectors, another property shared with DPMS.

The rest of this article is organized as follows. Section 2 presents our FV detection pipeline. Section 3 shows evidences that this model contains parts with a multimodal appearance. Section 4 gives insights on the mechanism that allows orderless statistics to capture this multimodality into a single model. Section 5 evaluates the level of sparsity contained in the model. Finally, Section 6 summarizes our findings.

2 Fisher-Vector detection

For the purpose of our analysis, we re-implemented the FV object detector described in [3]. To focus on the FV representation and avoid confounding factors, we did not employ the color descriptors, the contextual rescoring, or the local feature weighting using masks. Instead, we reproduced, and even slightly improved, the “baseline” version of their detector. We first describe this FV detection pipeline and then validate it empirically.

2.1 Detection pipeline

The *Fisher Vector* (FV), as used in computer vision applications [14], is obtained by aggregating the first and second order statistics of local descriptors, here SIFT [12], describing corresponding image patches. Given a Gaussian mixture model (GMM) with parameters μ_k (means), σ_k (diagonal covariance matrices), and π_k (priors), each D -dimensional SIFT descriptor x_i is first assigned to a mixture component k (following [3] we use hard assignments), and then the following first- and second-order statistics are computed:

$$\phi_{i,k}^{(1)} = \frac{1}{\sqrt{\pi_k}} \frac{x_i - \mu_k}{\sigma_k}, \quad \phi_{i,k}^{(2)} = \frac{1}{\sqrt{2\pi_k}} \left(\left(\frac{x_i - \mu_k}{\sigma_k} \right)^2 - 1 \right). \quad (1)$$

Every patch x_i is represented by the concatenation of the component statistics $\Phi_i = (\phi_{i,k}^{(1)}, \phi_{i,k}^{(2)} : k = 1, \dots, K)$ which we refer to as *point-wise FV* [2]. Its dimensionality is equal to $2KD$, where D is the dimension of x_i and K is the number of Gaussian components. Given an image region (*e.g.a* bounding box) from which we extract N patches, the final FV Φ is the average of point-wise FVs: $\Phi = \frac{1}{N} \sum_{i=1}^N \Phi_i$.

Weak geometry. To incorporate weak geometry in the representation we follow [3] and use a spatial pyramid [11]: each candidate image region is subdivided into 1×1 and $R \times R$ spatial subdivisions and the corresponding $R^2 + 1$ FVs are extracted and stacked.

Normalization. As shown in [14], the performance of the FV can be substantially improved by signed-square-rooting (SSR) each dimension followed by ℓ_2 normalization. The SSR is used in [3]. Here we have found that slightly better performance can be obtained by switching to *intra-normalization* [1, 15], *i.e.* by applying ℓ_2 normalization to each individual component statistics after aggregation:

$$\Phi' = \frac{1}{\sqrt{K}} \left[\frac{\Phi^1}{\|\Phi^1\|_2} \quad \dots \quad \frac{\Phi^K}{\|\Phi^K\|_2} \right], \quad \text{where} \quad \Phi^k = \left[\sum_{i=1}^N \phi_{(i,k)}^{(1)} \quad \sum_{i=1}^N \phi_{(i,k)}^{(2)} \right]. \quad (2)$$

Note that the $1/\sqrt{K}$ factor guarantees that $\|\Phi'\|_2 = 1$ when at least one patch is assigned to each Gaussian.

FV detection. So far, we have described the FV as a method to compute a descriptor for a particular image region. To use this to detect and localize

object occurrences, we could apply it in a sliding window manner over the image. For efficiency reasons, however, we follow [3] and, instead of trying all image subwindows, we limit the search to the ones enclosing the object region proposals obtained using the algorithm of [17].

2.2 Experimental setup

We evaluate our FV pipeline on the PASCAL VOC 2007 challenge dataset (VOC-07) [7], a standard test-bed that contains 20 categories of animals, vehicles and indoor objects. Following [3] we use a small subset of images from the VOC-07 training set – referred to as VOC-07-SMALL – involving only 4 classes to validate the parameters of our pipeline. This small set is also considered later for the most computationally-demanding experiments.

We extract 12×12 patches with a step size of three pixels at fifteen scales separated by a factor 1.2. SIFT [13] descriptors are extracted for each patch. We project and decorrelate these descriptors to $D = 64$ dimensions using PCA. The visual codebook consists of $K = 64$ Gaussian components and we use 1×1 and $R \times R$ non-overlapping spatial subdivisions with $R = 4$ in the spatial pyramid. This results in the concatenation of 17 intra-normalized FVs, yielding a $17 \times 2KD = 139K$ dimensional descriptor, which is then ℓ_2 -normalized again. We extract around 1500 candidate windows per image using selective search [17]. At train time, we apply 3 rounds of hard-negative mining extracting each time the top two false positive detections for each training image. At test time, non maximum suppression is applied in order to discard redundant detections that overlap more than 30% using the intersection-over-union overlap measure.

This pipeline obtains a mean average precision (or mAP) over the 20 classes of 33.9% with SSR, which is comparable to the 34.0% reported by [3], and a mAP of 34.6% with the *intra-normalization* scheme, thus confirming the positive effects of this normalization.

3 The Fisher Vector detector as a part-based model

A fundamental problem in visual recognition, and especially in object detection, is accounting for intra-class variability. This includes variations in appearance between two views of the same object instance (*e.g.* difference in view point), but also variations in appearance between two different instances of the same class (*e.g.* the class dog contains instances of Poodles and German Shepards). Consequently, images of the same object class form a *high-dimensional multimodal distribution* which is challenging to model.

As mentioned earlier, the DPM [8] improves on the basic Dalal-Triggs [5] HOG-detector in two key ways: by the introduction of *multiple aspects* and by the introduction of *deformable parts*. The use of multiple aspects addresses the multimodal issue: it allows the model to capture very different object appearances as caused for example by a large out-of-plane 3D rotation of the

object. The introduction of parts addresses the high-dimensional issue: the object is broken-down into smaller “pieces” which are easier to model because they typically lie in a lower dimensional space. The fact that the parts are deformable allows the template to warp geometrically and therefore adapt to image-based deformations of the object. This is important to ensure that the unimodal assumption of each aspect and low-dimensional assumption of each part are reasonable. To summarize, the **DPM is a mixture of aspects** or components, and each aspect is a **collection of parts**.

The DPM should be contrasted with the FV detector. The latter models the object appearance by extracting $1 + R^2$ FVs, one for each spatial subdivision in the pyramid. Each FV captures the appearance of a corresponding spatial bin by pooling local SIFT patches. If we define a part as one of these spatial bins, the model can be seen as a **collection of parts** where each part is modeled as a **mixture of Gaussians**. One could argue that individual SIFT patches could also be considered as parts, but Section 4 shows that these are significantly lower level, closer to part fragments.

In the following, we show that interpreting the FV detector as a part-based model explains how the FV can represent the high-dimensional multimodal appearance of object categories. We start by comparing parts in DPM and FV. Both models have a similar structure: one root part that captures information at the level of the whole object (the ‘root filter’ in DPM, and the 1×1 spatial bin of the pyramid corresponding to the whole object in FV) and several local parts (the ‘part filters’ in DPM, and the bins of the $R \times R$ spatial subdivisions in FV). Yet, their geometry is different. While parts in DPMs move in order to match the deforming structure of the underlying object, parts in the FV detector have a rigid predefined layout. A second key difference is that, because of the rigidity of the geometry and the lack of multiple components in the model, each part in the FV detector is required to capture a highly-variable and multimodal appearance.

The following experiment compares parts in the DPM and the FV detectors, illustrating the variability of the part appearances captured by the two models. Parts in the FV detector are visualized as follows: given a class – motorbike in our example – and the VOC-07-SMALL test set, the 200 top scoring object detections returned by the FV detector are selected. The images of these 200 detections as well as the subimages corresponding to the individual parts are extracted and, for each part individually, the corresponding FV descriptors are clustered in six groups using K-means. The average image of each group is shown in Figure 1. For comparison, a similar procedure is used to show the parts captured by a DPM model. The 200 top scoring detections are selected and split into six distinct sets according to the DPM component that was used for each detection. The images belonging to each group are then averaged and shown in Figure 2 for the whole object and for each part.

Given Figure 1 and Figure 2, we note that both models are capable of detecting with high confidence very dissimilar objects and object parts. However, in the FV detector multimodality is captured at the level of the individual parts, as shown by the highly-variable appearance of the clusters in Figure 1

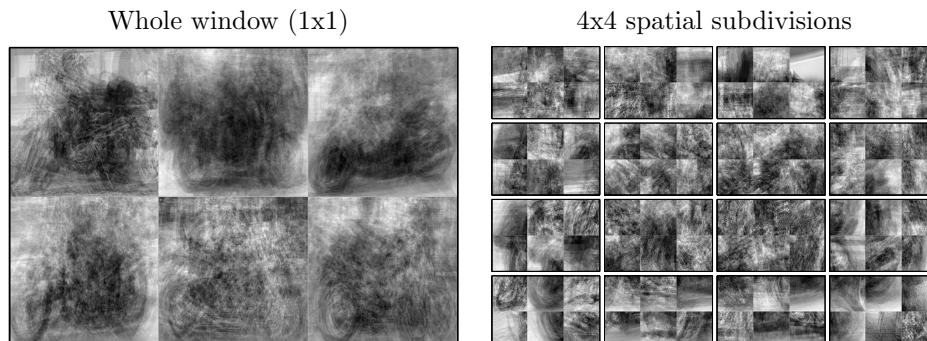


Figure 1: Left: average of the top 200 motorbike detections in VOC-2007-SMALL using the FV detector for the class ‘motorbike’. Six averages are shown corresponding to clustering the 200 images in six groups, as explained in the text. Right: the same procedure is applied to the individual parts. Note the significant variability of the part appearance.

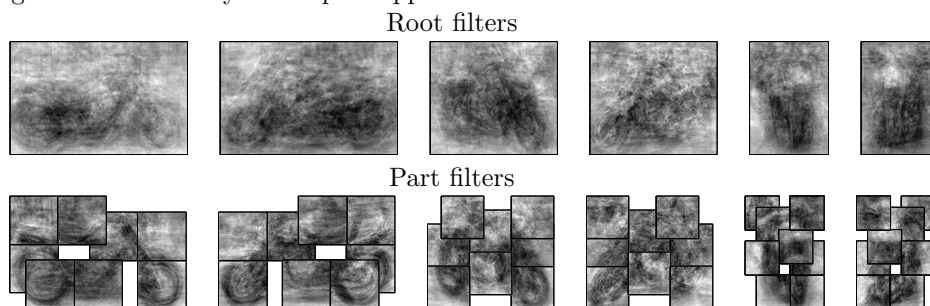


Figure 2: Average of the top 200 motorbike detections in VOC-2007-SMALL using the DPM model. Six averages are shown, each computed from the subset of the 200 images that activated a particular DPM component.

(we also tested averaging images without clustering, but this, as expected, blurs any detail). For the DPM, multimodality is captured instead at the level of the components of the mixture, which correspond to meaningful modes in the appearance space.

A direct conclusion of this observation is that the FV detector captures the object variability by means of a collection of multimodal parts rather than by using a mixture model as the DPM. Indeed, if we assume that each part has M modes¹ and if we have N parts, then the FV detector can model M^N different appearances of the same object. In other words, it is possible to model an exponential number of appearances because the multimodality is factorized per part. On the contrary, the DPM can have at most as many modes as components, because it is built as a mixture model by construction.

¹For the simplicity of the analysis, we assume that we have the same number of modes for each part.

Another difference between the DPM and FV detectors lies in training. The DPM learns a different part-based model for each component, where components are roughly divided according to viewpoints. In practice, components are initialized based on aspect-ratio or another heuristic, and eventually assembled in a mixture model using latent variables. In contrast, the FV detector learns in one shot a single linear classifier capturing the whole space of object appearances; nevertheless, the factorization property of the FV allows this procedure to capture efficiently an exponential number of different appearances.

4 Modeling multimodal parts with the Fisher Vector

In the previous section we have interpreted the FV detector as a part-based model, and we have given evidence that these parts are highly multimodal, a necessary property for dealing with the highly variable appearance of object categories. In this section we investigate the mechanism that allows FV to capture such rich multimodal appearances. In BoVW models, multimodality is easily explained as the representation quantizes the feature spaces in thousands of different visual words. However, the quantization granularity is much smaller for FV, typically in the order of a few dozens Gaussian components. This section clarifies why such a small number of visual words is sufficient to represent rich appearance variations.

The first answer to this question lies in the statistics encoded by the FV. While both BoVW and FV quantize local patch descriptors using a visual codebook, BoVW captures only 0th-order statistics (counts) of the features, while FV captures first order and second order statistics as well (see Eq. (1)). In particular, the point-wise FV $\phi(x_i)$ is a mixture of quadratic functions of the descriptor x_i , where different functions are activated based on the descriptor quantization. Hence, a linear classifier $\langle \mathbf{w}, \Phi(x_i) \rangle$ learnt on the FV representation can be seen as a mixture of quadratic experts in the space of local patch descriptors.

Our next experiment demonstrates the expressive power of the FV representation by showing that a linear classifier learned on top of FV can induce complex decision boundaries in descriptor space. To this end, we consider a patch x_i and encode it using the point-wise FV $\Phi(x_i)$. We then associate this patch x_i to a score $s(x_i|c, b) = \langle \mathbf{w}_{cb}, \frac{\Phi(x_i)}{\|\Phi(x_i)\|_2} \rangle$ where \mathbf{w}_{cb} denotes the weight vector learned by the FV detector for class c and spatial bin b . For the purpose of this visualization, we score a single patch at a time although the weights \mathbf{w}_{cb} are trained on the full FV model that pools information from all the patches covering an object. While in this manner we cannot visualize the aggregated effect of all the patches, Sect. 5 shows that only a small number of those is actually important for classification making this a good proxy.

Our goal is to illustrate the richness of the scoring functions $s(x_i|c, b)$ that the FV representation associates to local descriptors even when a single Gaussian

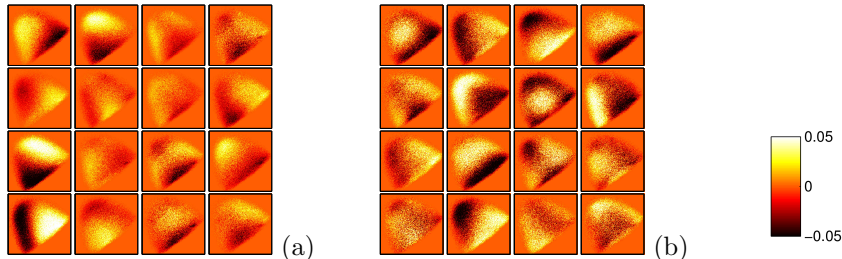


Figure 3: The FV can be seen as a mixture of quadratic experts in the space of local descriptors. The figure illustrates the variety of scoring functions that a single expert, corresponding to a single Gaussian in the model, can express. Each image shows the scoring functions obtained for each of the 4×4 spatial bins of (a) buses and (b) motorbikes.

mode is considered. To this end, we restrict the function domain to the patches $x_i \in X_k$ that are assigned to a certain Gaussian k . We then plot the function $s(x|c, b) : x \in X_k$ for different categories c and spatial bins b . Since the input of this function are 128 dimensional SIFT descriptors ($X_k \subset \mathbb{R}^{128}$), we first parametrize the descriptor space with a 2D space $\hat{x} = M_k x$, $\hat{x} \in \mathbb{R}^2$ where M_k is a PCA projection learnt from a set of descriptors sampled from X_k . Figure 3 shows $s(M_k x|c, b)$ for a Gaussian k , two classes c (bus and motorbike), and $4 \times 4 = 16$ spatial bins. These results are representative of other classes and Gaussian components.

There are two points to take from the visualizations in Figure 3. First, scores are well clustered in a small number of modes, due to the smooth form of the scoring function and of the encoding function $\Phi(\mathbf{x}_i)$. Second, the modes are nevertheless very varied, both for different classes and for different spatial bins, showing that the same Gaussian cluster has significantly different meaning depending on the class as well as on the spatial location. The ability of “reusing” the same visual words to express varied decision functions explains how the FV is able to capture complex multi-modal object appearances while using visual vocabularies significantly smaller than BoVW.

While Figure 3 shows that the FV is capable of assigning complex scoring functions to local patches, it does not clarify whether the learned scores are “semantic”, in the sense of identifying interpretable image fragments. Our next goal is to investigate this question. To do so we consider the set X_k of all image patches extracted from the VOC-07-SMALL dataset and assigned to the k^{th} Gaussian component of the model. Within this set, we select the 36 patches x_i that maximize the score $s(x_i|c, b)$ for a given category c and spatial bin b . We visualize these patches in Figure 4 for different objects and spatial bins. Despite the fact that all the patches in Figure 4 happen to belong to true positive detection windows, it is very hard to recognize which object parts they belong to. Hence, we conclude that the local patches pooled by the IFV are akin to *part fragments* rather than semantic parts; on the other hand, Figure 1 suggests that, by aggregating many of these fragments together, the spatial bins

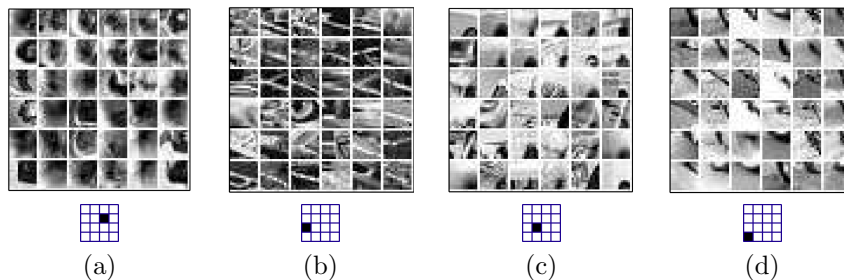


Figure 4: Highest scoring patches for some Gaussians and some spatial bins (bottom row) of the 4×4 spatial layout, for the cat (a), sheep (b), bus (c) and motorbike (d) classes.

in the FV can recognize meaningful parts.

We conclude that the FV captures multimodal part appearances by (1) decomposing parts into distributions of lower-level image fragments and (2) by learning complex scoring functions for these fragments despite the use of a small visual vocabulary.

5 Sparsity properties of the Fisher Vector detector

A desirable property of the DPM – at least on a visualization standpoint, but to some extent also on a computational point of view – is that it is sparse, in the sense that it captures the appearance of objects using a small set of part templates that fire in correspondence of selected image regions that are aligned to the parts (a consequence of the use of max-pooling). On a first glance, the FV detector does not seem to exhibit any sparsity property: patches are quantized using a GMM with a number of components that, while much smaller than visual words in BoVW, is still larger than the number of parts in a DPM; furthermore, thousands of image patches are encoded and averaged due to the use of sum-pooling instead of max-pooling. We revisit this assumption and study the sparsity properties of the FV detector at two different levels, the patch level and the Gaussian level.

Sparsity of the patches. Sect. 4 and Figure 3 looked at the scores associated by point-wise Fisher Vectors to individual top-ranked image patches. However, the FV detector pools information from hundreds of patches in the detection window and it is unclear whether the final decision depends on these relatively rare highly-scoring patches or, instead, the majority of other patches that receive intermediate scores. In other words, we do not know whether information is concentrated as it happens in the DPM case or instead distributed uniformly in the detection window. The next experiment answers this question.

We start from the hypothesis that top-ranked patches (*i.e.* those at the centers of the positive and negative modes in Figure 3) contain most of the infor-

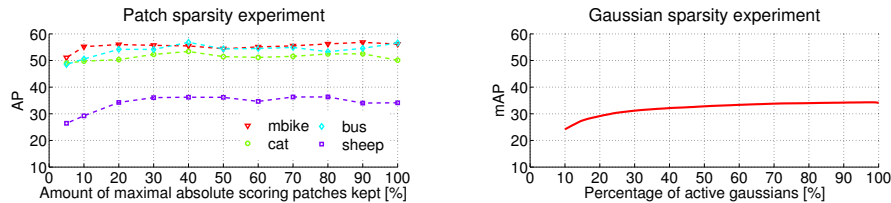


Figure 5: Sparsity experiments. Left: as an increasing number of low-scoring patches is removed from the model, detection AP changes very little until 80% of the patches are discarded (the AP of four representative classes is shown). Right: detection mAP (across all classes) remains stable by removing up to 50% of the Gaussian components from the model. See text for details.

mation necessary for detection. In order to validate this hypothesis, we drop from each detection window a certain portion of patches, starting from the ones with smallest absolute individual score, while keeping track of the achieved APs. Intuitively, if these low-scoring patches have a small effect on the final classifier score, the detection APs should be stable. However, once the important patches start to be removed APs should decrease rapidly.

The experiment was carried on the development set VOC-07-SMALL. The left panel in Figure 5 shows the results. It is apparent that even after removing about 80% of the patches with low absolute values of their scores the detector performance remains largely unchanged, thus confirming our intuition that only a small subset of patches is actually important for the classifier. We call this property *patch-level sparsity*.

Sparsity of the Gaussians components. We have shown earlier that the FV is capable of capturing complex decision boundaries despite making use of a small visual codebook. Here we investigate whether any of these visual words are in fact not informative and can be discarded during detection yielding, among other things, a computational saving. To do so we take advantage of the block structure of the FV. As explained in Section 2, each FV contains K non-overlapping blocks of size $2D$ corresponding to individual Gaussian components. The idea is then to induce sparsity at the block level when training the model, thereby encouraging the classifier to discard some of these Gaussian statistics.

In order to induce block sparsity, we learn the FV model vector \mathbf{w} using an SVM but replacing the standard ℓ_2 regularizer $\lambda \|\mathbf{w}\|_2^2$ with the *group lasso* one $\Omega(\mathbf{w})_{\ell_1/\ell_2}$ [21]. More precisely $\Omega_{\ell_1/\ell_2}(\mathbf{w}) = \lambda \sum_{g=1}^K \|\mathbf{w}_g\|_2$, where \mathbf{w}_g is the set of weights corresponding to Gaussian g in the FV statistics. Increasing the regularization strength λ favors models for which many Gaussians g have $\|\mathbf{w}_g\|_2 = 0$, effectively removing more Gaussians from the representation. The SVM objective function is optimized using the dual averaging method of [20]. A naive application of this method, however, results in a substantial performance drop as the sparsity of w increases. A possible reason for this loss of performance is that group lasso is capable of selecting useful component, but that the ℓ_1/ℓ_2 regularization is just not competitive with vanilla ℓ_2 . In order to validate

this hypothesis, we employ a fine-tuning step in which first group lasso is used to select a subset of useful Gaussians, and then a standard SVM classifier with ℓ_2 -norm regularization is used to learn the final detection model. Perhaps surprisingly, ℓ_2 retraining recovers most of the lost performance, and is therefore essential to obtain a good sparse model. For example, from the results of experiments on the development set (VOC-07-SMALL) we observed that when using the model which discards 50% of the Gaussians, there is a decrease of 8.3 mAP (*i.e.* the mean AP computed over all four classes from the development set) if the fine-tuning step is omitted. We think that these observations may transfer to several other applications of group lasso.

Overall Figure 5 shows that we can remove up to 50% of the Gaussians and still obtain comparable results to the full model. Note that here Gaussians are counted on a per-spatial-bin basis, as they are reused in different ones. As such, the 100% mark on x -axis of Figure 5 corresponds to $(4 \times 4 + 1) \times 64 = 1088$ active Gaussians. Nevertheless, eliminating Gaussian components allows us to avoid accumulating corresponding point-wise FVs during the detection phase yielding a proportional acceleration in detection (note that patches can be quantized once for all detection windows in an image, but accumulation occurs for each candidate window separately).

6 Summary

In this paper, we have shown that the FV detector contains parts in the same manner as a DPM. In both cases, a fixed number of parts is used to capture the diversity of appearance of an object category. However, this diversity is represented differently: the DPM uses a mixture of components where each component corresponds to an aspect, and the FV factors the appearance in a product of multimodal distributions, one for each object part, encoding implicitly an exponential number of combinations. Both DPM and FV are sparse, although in somewhat different senses. DPM is sparse at the level of the parts, which are few and max-pooled, while FV exhibits sparsity at the level of part fragments and at the level of the visual vocabulary. For example, while FV has typically more parts than DPM, 80% of the pooled patches and 50% of the visual words can be removed with minimal impact on performance. The latter fact can be used in order to accelerate detection proportionally.

References

- [1] R. Arandjelovic and A. Zisserman. All about VLAD. In *Proc. CVPR*, 2013.
- [2] Q. Chen, Z. Song, R. Feris, A. Datta, L. Cao, Z. Huang, and S. Yan. Efficient maximum appearance search for large-scale object detection. In *CVPR*, 2013.

- [3] R. G. Cinbis, J. Verbeek, and C. Schmid. Segmentation driven object detection with fisher vectors. In *Proc. ICCV*, 2013.
- [4] G. Csurka, C. R. Dance, L. Dan, J. Willamowski, and C. Bray. Visual categorization with bags of keypoints. In *Proc. ECCV Workshop on Stat. Learn. in Comp. Vision*, 2004.
- [5] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Proc. CVPR*, 2005.
- [6] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2009 (VOC2009) Results. <http://www.pascal-network.org/challenges/VOC/voc2009/workshop/index.html>, 2009.
- [7] M. Everingham, A. Zisserman, C. Williams, and L. Van Gool. The PASCAL visual object classes challenge 2007 (VOC2007) results. Technical report, Pascal Challenge, 2007.
- [8] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *PAMI*, 2009.
- [9] P. F. Felzenszwalb, D. McAllester, and D. Ramanan. A discriminatively trained, multiscale, deformable part model. In *Proc. CVPR*, 2008.
- [10] H. Jégou, M. Douze, C. Schmid, and P. Pérez. Aggregating local descriptors into a compact image representation. In *Proc. CVPR*, 2010.
- [11] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bag of features: Spatial pyramid matching for recognizing natural scene categories. In *Proc. CVPR*, 2006.
- [12] D. G. Lowe. Object recognition from local scale-invariant features. In *Proc. ICCV*, 1999.
- [13] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 2(60):91–110, 2004.
- [14] F. Perronnin, J. Sánchez, and T. Mensink. Improving the Fisher kernel for large-scale image classification. In *Proc. ECCV*, 2010.
- [15] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep fisher networks for large-scale image classification. In *Proc. NIPS*, 2013.
- [16] J. Sivic and A. Zisserman. Video Google: A text retrieval approach to object matching in videos. In *Proc. ICCV*, 2003.
- [17] J.R.R. Uijlings, K.E.A. van de Sande, T. Gevers, and A.W.M. Smeulders. Selective search for object recognition. *International Journal of Computer Vision*, 2013.

- [18] A. Vedaldi, V. Gulshan, M. Varma, and A. Zisserman. Multiple kernels for object detection. In *Proc. ICCV*, 2009.
- [19] C. Vondrick, A. Khosla, T. Malisiewicz, and A. Torralba. HOGgles: Visualizing object detection features. In *Proc. ICCV*, 2013.
- [20] H. Yang, Z. Xu, I. King, and M. R. Lyu. Online learning for group lasso. In *IMCL*, 2010.
- [21] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.