

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

# Fisher vector places: learning compact descriptors for place recognition

Anonymous CVPR submission

Paper ID 954

## Abstract

The goal of this work is to localize a query photograph by finding other images depicting the same place in a large geotagged image database. The contributions of this work are three fold. First, we learn a discriminative yet compact descriptor of each image in the database. This is achieved by applying exemplar support vector machine (e-SVM) learning to compact Fisher vector descriptors extracted from database images. Secondly, we analyze the exemplar support vector machine objective and show that the learnt hyperplane can be interpreted as a new descriptor that replaces the original positive example and is re-weighted to increase its separation from the negative data. Finally, based on this analysis we demonstrate that the learnt and re-normalized descriptor could be directly used for matching, thus avoiding the need for expensive and tedious calibration typically needed for exemplar support vector machine methods. Place recognition results are shown on two image datasets of Google street-view images from Pittsburgh and demonstrate the learnt representation consistently improves over the standard Fisher vector descriptors at different target dimensions.

## 1. Introduction

The goal of this work is to localize a query image by matching to a large database of geotagged street-level imagery. This is an important problem with practical applications in robotics, augmented reality or navigation. This task is however very difficult. It is hard to distinguish different places, e.g. streets in a city, from each other. The imaged appearance of a place can change drastically due to factors such as viewpoint, illumination or even changes over time. Finally, with the emergence of planet-scale geotagged image collections, such as Google Street-view, the image databases are becoming very large. We estimate a single country like France is covered by more than 60 million street-level panoramas. Hence the fundamental challenge in place recognition lies now in designing robust, discriminative, yet compact, image representations.

In this work we build on the method of Gronat *et al.* [12] who represent each image in the database by a per-location classifier that is trained to discriminate each place from other places in the database. At query time, the query image is classified by all per-location classifiers and assigned to a place with the highest classification score. The training of each classifier is performed using the per-exemplar support vector machine (e-SVM) [22], which takes the positive image as a single positive example and other far away images in the database as negative examples. The exemplar SVM is well suited for this task as street-level image collections typically contain only one or at most a hand-full of images depicting the same place. The intuition is that the exemplar SVM can learn the important features that distinguish the particular place from other similar places in the database. While the results of [12] are promising they suffer from two important drawbacks. First, the learnt place specific representation is not compact, which prohibits its application to planet-scale street-level collections that are now becoming available [16]. Second, the per-exemplar classifiers require careful and time-consuming calibration.

In this work we address both these issues. First, we apply the exemplar SVM training to compact Fisher vector [14, 24] image descriptors, which results in a *discriminative* yet *compact* representation of each image in the database. Second, to avoid the expensive classifier calibration, we analyze the exemplar SVM cost and show that the learnt hyperplane can be interpreted as a new descriptor that replaces the original positive example and is re-weighted to increase its separation from the negative data. As a result of this analysis, we demonstrate that after an appropriate normalization of the new re-weighted descriptor no further calibration is necessary. We show improved results on place recognition using learnt compact Fisher vector descriptors [14] of different dimensionality.

## 2. Related work

**Large-scale visual place recognition.** The visual localization problem is typically treated as a large-scale instance-level retrieval [6, 4, 12, 17, 27, 34, 35], where images are represented using local invariant features [21] aggre-

108 gated into the bag-of-visual-words [5, 32] representation.  
 109 The image database can be further augmented by 3D point  
 110 clouds [16], automatically reconstructed by large-scale  
 111 structure from motion (SfM) [1, 16], which enables accurate  
 112 prediction of query image camera position [20, 25].  
 113 In this work we investigate learning a discriminative  
 114 representation using the compact Fisher vector descriptors [15].  
 115 Fisher vector descriptors have shown excellent place recog-  
 116 nition accuracy [34]. In this work we further improve their  
 117 performance by discriminative learning.  
 118

119  
 120 **Fisher vector image representations.** Fisher vector im-  
 121 age representations have recently demonstrated excellent  
 122 performance for a number of visual recognition tasks [3,  
 123 15, 18, 30]. They are specially suited for retrieval applica-  
 124 tions since they are robust to image appearance variations  
 125 and capture richer image statistics than the simple bag-of-  
 126 visual-words (BOW) aggregation. However, the raw ex-  
 127 tracted Fisher vectors are typically high-dimensional, e.g.  
 128 with 32,768 non-sparse dimensions, which is impractical  
 129 for large-scale visual recognition and indexing applications.  
 130 Hence, their dimensionality is often reduced by principal  
 131 component analysis (PCA) and further quantized for effi-  
 132 cient indexing using, e.g. a product quantizer [15]. Other  
 133 recent work has demonstrated improved performance in a face  
 134 recognition application by finding discriminative projection  
 135 using a large amount of training face data [30]. Our work  
 136 is complementary to these methods as it operates on the  
 137 projected low-dimensional descriptor and further learns dis-  
 138 criminative re-weighting of the descriptor specific to each  
 139 image in the database using per-exemplar support vector  
 140 machine [22].  
 141

142  
 143 **Per-exemplar support vector machine.** The exemplar  
 144 support vector machine (e-SVM) has been used in a number  
 145 of visual recognition tasks including category-level recog-  
 146 nition [22], cross-domain retrieval [29], scene parsing [33],  
 147 place recognition [12] or as an initialization for more com-  
 148 plex discriminative clustering models [8, 31]. The main  
 149 idea is to train a linear support vector machine (SVM) clas-  
 150 sifier from a single positive example and a large number of  
 151 negatives. The intuition is that the resulting weight vector  
 152 will give a higher weight to the discriminative dimensions  
 153 of the positive training data point and will down weight  
 154 dimensions that are non-discriminative with respect to the  
 155 negative training data. A key advantage is that each per-  
 156 exemplar classifier can be trained independently and hence  
 157 the learning can be heavily parallelized. The per-exemplar  
 158 training brings however also an important drawback. As  
 159 each classifier is trained independently a careful calibration  
 160 of the resulting classifier scores is required [12, 22].  
 161

162 **Contributions.** The contributions of this work are three-  
 163 fold. First, we analyze the exemplar support vector machine  
 164 objective and show that the learnt hyperplane can be inter-  
 165 preted as a new descriptor that replaces the original positive  
 166 example and is re-weighted to increase its separation from  
 167 the negative data. Secondly, we demonstrate that after an  
 168 appropriate normalization of the new re-weighted descriptor  
 169 no further calibration is necessary. Finally, we apply  
 170 e-SVM training to compact Fisher vector descriptors for  
 171 large-scale place recognition resulting in a *discriminative*  
 172 yet *compact* representation of each image in the database.  
 173 Place recognition results are shown on a dataset of 25k im-  
 174 ages of Pittsburgh and demonstrate the learnt representation  
 175 consistently improves over the standard Fisher vector de-  
 176 scriptors at different target dimensions.  
 177

### 3. Learning compact place descriptors using per-exemplar SVM

178 Each database image  $j$  is represented by its L2-  
 179 normalized Fisher vector  $\Phi_j$ . The goal is to learn a set  
 180 of new L2-normalized Fisher vectors  $\Psi_j$ , one per each  
 181 database image. At query time, given the Fisher vector  $\Phi_q$   
 182 of an unknown query image, we retrieve the database image  
 183 depicting the same location by finding the image  $j^*$  with the  
 184 highest score measured by a dot product  
 185

$$j^* = \arg \max_j \Phi_q^T \Psi_j. \quad (1)$$

186 In other words, the aim is to replace each original database  
 187 Fisher vector  $\Phi_j$  with a new vector  $\Psi_j$  that is more dis-  
 188 criminative in the sense of separation from descriptors of  
 189 images depicting other places. Inspired by [12], we in-  
 190 vestigate applying the exemplar support vector machine (e-  
 191 SVM) [22] for this task. e-SVM learns a linear classifier  
 192  $w_j^T \Phi + b_j$  given the descriptor  $\Phi_j^+$  of place  $j$  as a single  
 193 positive example (with target label +1) and a large number  
 194 of negative descriptors  $\mathcal{N}_j$  from other places in the database  
 195 (with target labels -1). The intuition of the exemplar SVM  
 196 training [22] is that the learnt weight vector  $w_j$  will give a  
 197 higher weight to the dimensions of the descriptor that are  
 198 discriminative and will down-weight dimensions that are  
 199 non-discriminative with respect to the negative training data  
 200 collected from other far-away places. The optimal  $w_j$  and  
 201  $b_j$  are obtained by minimizing the following objective  
 202

$$\begin{aligned} & \|w_j\|^2 + C_1 \cdot h(w_j^T \Phi_j^+ + b_j) \\ & + C_2 \sum_{\Phi \in \mathcal{N}_j} h(-w_j^T \Phi - b_j), \end{aligned} \quad (2)$$

203 where  $\Phi_j^+$  is the descriptor of the place  $j$  as the positive data  
 204 point,  $\Phi$  are Fisher descriptors from negative training data  
 205  $\mathcal{N}_j$  and  $h$  is the hinge loss,  $h(y) = \max(1 - y, 0)$ . Note  
 206 that the first term in (2) is the regularizer, the second term  
 207

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
is the loss on the positive data weighted by scalar parameter  $C_1$  and the third term is the loss on the negative data weighted by scalar parameter  $C_2$ . The objective is convex and can be minimized with respect to  $w_j$  and  $b_j$  using standard software packages such as [10]. A key advantage is that the per-exemplar classifier for each place can be trained independently and hence the learning can be heavily parallelized. The downside of the independent training for each positive example is that the resulting scores have to be calibrated with respect to each other on additional data [12, 22].

### 3.1. Analysis of per-exemplar SVM objective

In this section, we analyze the exemplar SVM objective (2) and show the learnt and *re-normalized* weight vector  $w_j$  can be interpreted as a new descriptor  $\Psi_j$  that replaces the original positive training descriptor  $\Phi_j^+$ . In particular, we show first that when the weight  $C_2$  of the negative data in objective (2) goes to zero and the learnt  $\Psi_j$  is identical to the original positive training data point  $\Phi_j^+$ . Second, when  $C_2 > 0$ , the learnt  $\Psi_j$  moves away from the positive  $\Phi_j^+$  to increase its separation from the negative data. Details are given next.

**Case I:**  $C_2 \rightarrow 0$ . The goal is to show that when the weight  $C_2$  of the negative data in objective (2) goes towards zero, the resulting hyperplane vector  $w_j$  is parallel with the vector of positive training descriptor  $\Phi_j^+$ . When  $w_j$  is normalized to have unit L2 norm the two vectors are identical. First, let us decompose  $w$  into parallel and orthogonal part with respect to the positive training data point  $\Phi^+$  (in the following we omit index  $j$  for brevity), i.e.  $w = w^\perp + w^\parallel$ , where  $(w^\perp)^T \Phi^+ = 0$ . Next, we observe that when the weight of the negative data diminishes ( $C_2 \rightarrow 0$ ), any non-zero component  $w^\perp$  will increase the value of the objective. As a result, for  $C_2 \rightarrow 0$  the objective is minimized by  $w^\parallel$ , i.e. the optimal  $w$  is parallel with  $\Phi^+$ .

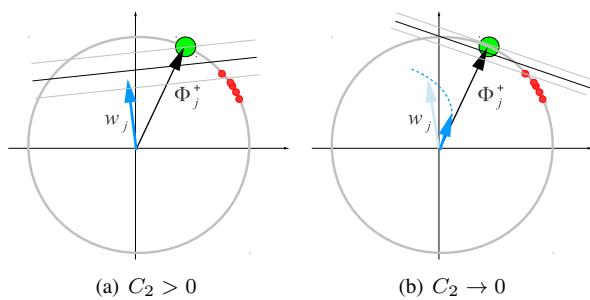
In detail, for  $w = w^\perp + w^\parallel$  the objective (2) can be written as

$$\begin{aligned} & \|w^\perp + w^\parallel\|^2 + C_1 \cdot h((w^\perp + w^\parallel)^T \Phi^+ + b) \\ & + C_2 \sum_{\Phi \in \mathcal{N}} h(-(w^\perp + w^\parallel)^T \Phi - b). \end{aligned} \quad (3)$$

Note that the orthogonal part  $w^\perp$  does not change the value of the second term in (3) because  $(w^\perp + w^\parallel)^T \Phi^+ = (w^\parallel)^T \Phi^+$ , and hence (3) reduces to

$$\begin{aligned} & \|w^\perp + w^\parallel\|^2 + C_1 \cdot h(w^\parallel T \Phi^+ + b_j) \\ & + C_2 \sum_{\Phi \in \mathcal{N}} h(-(w^\perp + w^\parallel)^T \Phi - b). \end{aligned} \quad (4)$$

In the limit case as  $C_2 \rightarrow 0$  any non-zero component  $w^\perp$  will increase the value of the objective (4). This can be



270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
Figure 1: An illustration of the effect of decreasing parameter  $C_2$  in the exemplar support vector machine objective. The positive exemplar  $\Phi_j^+$  is shown in green. The negative data points are shown in red. All training data is L2 normalized to lie on a hyper-sphere. (a) For  $C_2 > 0$ , the normal  $w_j$  of the optimal hyper-plane moves away from the direction given by the positive example  $\Phi_j^+$  in a manner that reduces the loss on the negative data. (b) As the parameter  $C_2$  decreases the learnt  $w_j$  becomes parallel to the positive training example  $\Phi_j^+$  and its magnitude  $\|w_j\|$  goes to 0.

seen by noting that the third term vanishes when  $C_2 \rightarrow 0$  and hence the objective is dominated by the first two terms. Further, the second term in (4) is independent of  $w^\perp$ . Finally, the first term will always increase for any non-zero value of  $w^\perp$  as  $\|w^\perp + w^\parallel\|^2 \geq \|w^\parallel\|^2$  for any  $w^\perp \neq 0$ .

As a result, in the limit case when  $C_2 \rightarrow 0$  the optimal  $w$  is parallel with  $\Phi^+$ . Note also, that when  $C_2$  is exactly equal to zero,  $C_2 = 0$ , the optimal  $w$  vanishes, i.e. the objective (4) is minimized by trivial solution  $\|w\| = 0$  and  $b_j = -1$ . The effect of decreasing the parameter  $C_2$  is illustrated in figure 1.

**Case II:**  $C_2 > 0$ . When the weight  $C_2$  of the negative data in the objective (4) increases the direction of the optimal  $w$  will be different from  $w^\parallel$  and will change to take into account the loss on the negative data points. Explicitly writing the hinge-loss  $h(x) = \max(1 - x, 0)$  in the last term of (4), we see that  $w$  will move in the direction that reduces  $\sum_{\Phi \in \mathcal{N}} \max(1 + w^\parallel T \Phi + b, 0)$ , i.e. that reduces the dot product  $w^\parallel T \Phi$  on the negative examples that are active (support vectors).

### 3.2. Interpreting normalized $w$ as a new descriptor

The above analysis demonstrates that as  $C_2$  decreases the normal of the optimal hyperplane  $w$  that separates the positive exemplar  $\Phi^+$  from negative data becomes parallel with  $\Phi^+$ , as shown in figure 1. As  $C_2$  increases, the normal  $w$  of the optimal hyper-plane moves away from the direction given by the positive example  $\Phi^+$  in a manner that reduces the loss on the negative data. This suggests that the learnt  $w$

324 could be interpreted as a modified positive example  $\Phi^+$ , re-  
 325 weighted to emphasize directions that separate  $\Phi^+$  from the  
 326 negative data. As discussed above  $w$  is not normalized. As  
 327 we wish to measure the similarity between descriptors by  
 328 (the cosine of) their angle given by equation (1), additional  
 329 normalization of the learnt  $w$  is necessary. Hence we define  
 330 the new descriptor  $\Psi_j$  as the normalized hyperplane normal  
 331  $w_j$

$$\Psi_j = \frac{w_j}{\|w_j\|}. \quad (5)$$

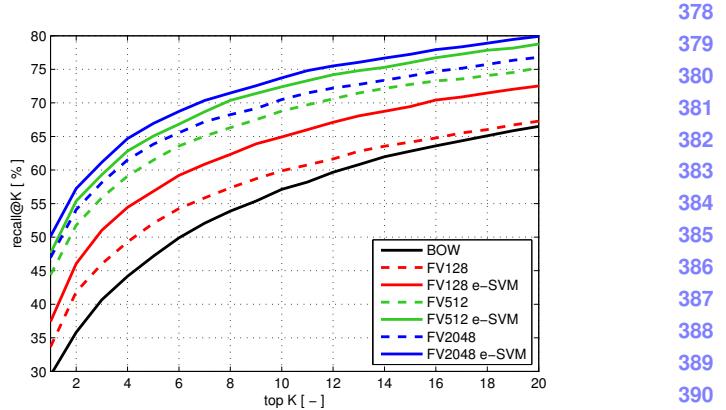
335 The proposed re-normalization procedure can be interpreted  
 336 as a simple form of affine calibration where the un-  
 337 calibrated SVM score  $s(x) = w^T x + b$  is transformed into a  
 338 calibrated score  $f(x) = \alpha s(x) + \beta$ , where  $\alpha = 1/\|w\|$  and  
 339  $\beta = -b\alpha$ . Note that the bias term  $b$  is not ignored but used  
 340 to compute  $f(x)$ . The intuition is that for L2 normalized de-  
 341 scriptors  $x$  computing the calibrated score  $f(x)$  corresponds  
 342 to computing the normalized cross-correlation between vec-  
 343 tors  $w$  and  $x$ . This was found to work well in retrieval (see  
 344 e.g. Sivic and Zisserman, ICCV 2003 [32]) or matching  
 345 whitened HOG descriptors [11,13] (see e.g. [Doersch et  
 346 al., Mid-Level Visual Element Discovery as discriminative  
 347 Mode Seeking, NIPS 2013]), but has not yet been used in  
 348 conjunction with SVMs as done in this work. Note also that  
 349 the affine calibration can be computed offline and does not  
 350 need any additional computation or storage at query time  
 351 compared to the p-value calibration of Gronat et al. [12].

## 353 4. Experimental evaluation

355 In this section we first describe the experimental set-up,  
 356 then we give implementation details, and finally report re-  
 357 sults of the proposed approach on two datasets where we  
 358 compare performance with raw Fisher vector matching and  
 359 several baselines methods.

### 360 4.1. Experimental set-up

362 We perform experiments on a database of Google Street  
 363 View images of Pittsburgh downloaded from the Internet.  
 364 The data contains panoramas covering roughly an area of  
 365  $1.3 \times 1.2 \text{ km}^2$ . Similar to [4], for each panorama we  
 366 generate 12 overlapping perspective views corresponding  
 367 to two different elevation angles to capture both the street-  
 368 level scene and the building façades, resulting in a total of  
 369 24 perspective views each with  $90^\circ$  FOV and resolution of  
 370  $960 \times 720$  pixels. For evaluation we have used two versions  
 371 of this data. The first one was obtained from the authors  
 372 of [12] (25k images). In the second version we download  
 373 additional images to increase the dataset size to 55k images.  
 374 As a query set with known ground truth GPS positions, we  
 375 use the 8999 panoramas from the Google Street View re-  
 376 search dataset. This dataset covers approximately the same  
 377 area, but has been captured at a different time, and depicts



395 **Figure 2: Evaluation on Pittsburgh 25k [12] dataset.** The  
 396 fraction of correctly recognized queries (recall@K, y-axis)  
 397 vs. the number of top  $K$  retrieved database images for  
 398 different Fisher vector dimensions. The learnt descriptors  
 399 by the proposed method (FV e-SVM) consistently improve  
 400 over the raw Fisher vector descriptors across the whole  
 401 range of  $K$  and all dimensions.

402 the same places from different viewpoints and under differ-  
 403 ent illumination conditions. For each test panorama, we  
 404 generate a set of perspective images as described above. Fi-  
 405 nally, we randomly select out of all generated perspective  
 406 views a subset of 4k images, which is used as a test set to  
 407 evaluate the performance of the proposed approach. Since  
 408 all the query images have associated GPS location we can  
 409 compute their spatial distance from the database images re-  
 410 turned by the matching method. We consider a query image  
 411 to be correctly localized if the retrieved database image lies  
 412 within a perimeter of 20m from the location of the query.

### 413 4.2. Implementation details

415 We first extract rootSIFT descriptors [2] for each im-  
 416 age. Following [15] we project the 128-dimensional SIFT  
 417 descriptors to 64 dimensions using PCA. The projection  
 418 matrix is learnt on a set of descriptors from 5,000 ran-  
 419 domly selected database images. This has also the effect  
 420 of decorrelating the SIFT descriptor. The 64-dimensional  
 421 SIFT descriptors are then aggregated into Fisher vectors us-  
 422 ing a Gaussian mixture model with  $N = 256$  components,  
 423 which results in a  $2 \times 256 \times 64 = 32,768$ -dimensional  
 424 descriptor for each image. The Gaussian mixture model  
 425 is learnt from descriptors extracted from 5,000 randomly  
 426 sampled database images. The high-dimensional Fisher  
 427 vector descriptors are then projected down to dimension  
 428  $d \in \{128, 512, 2048\}$  using PCA learnt from all avail-  
 429 able images in the database. The resulting low dimensional  
 430 Fisher vectors are then re-normalized to have unit L2-norm,  
 431 which we found to be important in practice.

**Learning parameters and training data.** To learn the exemplar support vector machine for each database image  $j$ , the positive and negative training data are constructed as follows. The *negative training set*  $\mathcal{N}_j$  is obtained by: (i) finding the set of images with geographical distance greater than 200m; (ii) sorting the images by decreasing value of similarity to image  $j$  measured by the dot product between their respective Fisher vectors; (iii) taking the top  $N = 500$  ranked images as the negative set. In other words, the negative training data consists of the hard negative images, i.e. those that are similar to image  $j$  but are far away from its geographical position, hence, cannot have the same visual content. The *positive training set*  $\mathcal{P}_j$  consist of the original Fisher vector  $\Phi_j$  of the image  $j$ . For the SVM training we use libsvm [10]. We use the same  $C_1$  and  $C_2$  parameters for all per-exemplar classifiers, but find the optimal value of the parameters for each dimensionality of the Fisher vector by a grid search evaluating performance on a held out set. We observe that for different Fisher vector target dimensions the optimal value of parameter  $C_1$  is quite stable (typically  $C_1 = 1$ ) while the optimal parameter for  $C_2$  varies between  $10^{-6}$  to  $10^{-1}$ . To learn the new image representation for each database image  $j$  we: (i) learn SVM from  $\mathcal{P}_j$  and  $\mathcal{N}_j$  (see above); (ii)  $L2$  normalize the learned  $w_j$  using equation (5); and (iii) use this re-normalized vector as the new image descriptor  $\Psi_j$  for image  $j$ . At query time we compute the Fisher vector  $\Phi_q$  of the query image and measure its similarity score to the learnt descriptors  $\Psi_j$  for each database image by equation (1).

### 4.3. Results

For each database (Pittsburgh 25k and Pittsburgh 55k) we compare results of our method (FV e-SVM) to two baselines: standard bag-of-visual-words baseline (BOW) and raw Fisher vector matching without learning (FV).

We perform experiments on several target Fisher vector dimensions  $d \in \{128, 512, 2048\}$ . For each method we measure performance using the percentage of correctly recognized queries (Recall) similarly to, e.g., [4, 17, 26]. The query is correctly localized if at least one of the top  $K$  retrieved database images is within 20 meters from the ground truth position of the query. Results are shown for different values of  $K$  in table 1. For the Pittsburgh 25k we also show results in the form of a curve in figure 2. The results clearly demonstrate the benefits of the learnt descriptors with respect to the standard Fisher vectors for all target dimensions and lengths of shortlist  $K$ . The benefits of discriminative learning are specially prominent for low-dimensional compact descriptors ( $d = 128$ ). The proposed method also significantly outperforms the bag-of-visual-words baseline. Figure 4 shows examples of place recognition results.

**Applying the proposed method to other descriptors** 486  
**TODO:** Applying the proposed method to other descriptors (R18, R38) 487  
The proposed approach can be applied 488  
to other descriptors beyond Fisher vectors. To demonstrate 489  
this we have applied the proposed calibration by re-normalization 490  
to the bag-of-visual-words descriptor used in 491  
(Gronat et al.[12]) and observed improvement from 29.4% 492  
recall@K=1 (vanilla bag-of-visual-words) to 32.6% recall@K=1 493  
(e-SVM with the proposed calibration by renormalization). 494  
For bag-of-visual-words, the e-SVM with p-value 495  
calibration gives comparable results (33.6% recall@K=1) 496  
but at an increased training cost and a significantly 497  
higher required memory footprint at query time, 498  
which makes the method very hard to scale-up to database 499  
sizes beyond 200k images. 500

**Comparison to other methods.** 501  
On the Pittsburgh 25k 502  
database, we compare performance of our learnt discriminative 503  
descriptors to the methods of [12] and [17], who report 504  
on the same testing data top  $K = 1$  recall of 36.5% and 505  
41.9%, respectively (results taken from [12]). Our method 506  
outperforms [17] already for dimension  $d = 128$  (37.8%) 507  
and [12] for dimension  $d = 512$  (47.6%). Furthermore, 508  
note that [12, 17] are based on a bag-of-visual-words 509  
representation, which typically needs to store between 1000-510  
2000 non-zero visual words per image, which is significantly 511  
more than our learnt 128 or 512 dimensional descriptor. 512

**TODO:** 513  
We have compared the proposed method with 514  
the p-value calibration of (Gronat et al.[12]) applied to 515  
Fisher vectors. The results show that our method significantly 516  
outperforms the p-value calibration (37.8% recall@1 517  
vs. 32.3% recall@1; 56.9% recall@5 vs. 50.0% recall@5) 518  
for the Fisher vector of dimension 128. Similar gains are 519  
observed for other FV dimensions. 520

**Memory complexity analysis.** 521  
Figure 3 compares the 522  
performance of the learnt discriminative descriptor (FV 523  
eSVM) to the raw Fisher vectors (FV) for different target 524  
dimensions. The results demonstrate that for a given level 525  
of accuracy (y-axis) our method learns a more compact 526  
(lower-dimensional) representation (x-axis). For example, 527  
our learnt 128-dimensional descriptor achieves a similar 528  
accuracy (around 65%) to the 256-dimensional raw Fisher 529  
descriptor essentially reducing the memory complexity to 530  
50% for the same level of performance. Note that similar 531  
to [15], we observe decrease in performance at high-532  
dimensions for both the baseline and our method. 533

### 5. Conclusions

We have shown that a discriminative yet compact image 536  
representation for place recognition can be learnt using the 537



Figure 4: **Examples of correctly and incorrectly localized queries for the learnt Fisher vector representation.** Each example shows a query image (left) together with correct (green) and incorrect (red) matches from the database obtained by the learnt Fisher vector representation *w-norm* method (top) and the standard Fisher vector baseline (bottom) for dimension 128. Note that the proposed method is able to recognize the place depicted in the query image despite changes in viewpoint, illumination and partial occlusion by other objects (trees, lamps) and buildings. Note that the baseline methods often finds images depicting the same buildings but in a distance whereas our learnt representation often finds a closer view better matching the content of the query.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

Method:	25k Pittsburgh					55k Pittsburgh				
	1	2	5	10	20	1	2	5	10	20
BOW	29.4	35.7	47.0	57.1	66.5	8.7	11.0	17.3	22.8	25.4
BOW raw SVM	6.4	8.1	13.5	17.5	20.5					
BOW e-SVM	32.2	37.3	46.2	54.4	63.3					
FV128	33.6	41.8	52.0	59.8	67.7	10.9	14.1	20.2	26.4	33.2
FV128 raw SVM	32.4	41.1	52.6	60.4	68.4					
<b>FV128 e-SVM</b>	<b>37.8</b>	<b>46.1</b>	<b>56.9</b>	<b>64.9</b>	<b>72.6</b>	<b>13.5</b>	<b>17.7</b>	<b>25.0</b>	<b>31.8</b>	<b>39.0</b>
FV512	44.3	51.7	61.4	68.7	75.2	17.3	21.1	28.4	34.2	40.3
<b>FV512 e-SVM</b>	<b>47.6</b>	<b>55.4</b>	<b>65.1</b>	<b>72.4</b>	<b>78.8</b>	<b>19.8</b>	<b>25.1</b>	<b>32.7</b>	<b>38.7</b>	<b>46.0</b>
FV2048	46.9	54.1	63.8	70.5	76.8	19.2	23.5	29.9	35.2	41.9
<b>FV2048 e-SVM</b>	<b>50.2</b>	<b>57.3</b>	<b>67.0</b>	<b>73.8</b>	<b>78.0</b>	<b>20.8</b>	<b>25.9</b>	<b>33.1</b>	<b>38.7</b>	<b>45.9</b>

Table 1: The fraction of correctly recognized queries (recall@K) vs. the number of top  $K \in \{1, 2, 5, 10, 20\}$  retrieved database images for different Fisher vector dimensions  $d \in \{128, 512, 2048\}$ . The learnt descriptors by the proposed method (FV e-SVM) consistently improve over the raw Fisher vector descriptors across the whole range of  $K$  and all dimensions on both the 25k and 55k Pittsburgh image datasets.

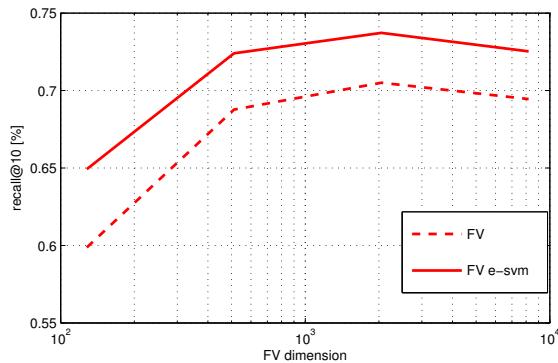


Figure 3: **Memory complexity analysis.** The fraction of correctly localized queries at the top 10 retrieved images (y-axis) for different Fisher vector dimensions (x-axis). The learnt descriptors by our method (FV e-SVM) clearly outperform the raw Fisher vector descriptors (FV) for all dimensions. Note that for a certain level of performance (y-axis) the proposed method learns a more memory efficient (lower dimensional, x-axis) descriptor.

exemplar support vector machine applied to Fisher vector image descriptors without the need for expensive and tedious calibration typical for other exemplar support vector machine methods. We demonstrate that proposed method is applicable to two different descriptors, the bag-of-visual-words and Fisher vectors. Our results show significant gains in place recognition performance compared to raw Fisher vector matching as well as other baselines. Our work opens up the possibility of learning a compact and discriminative representation using other descriptors such as HOG [7] or the recently developed convolutional neural network features [9, 19, 23, 28] as well as extending the analysis to

other cost functions [11, 13].

## References

- [1] S. Agarwal, N. Snavely, I. Simon, S. Seitz, and R. Szeliski. Building Rome in a day. In *ICCV*, pages 72–79, 2009. [2](#)
- [2] R. Arandjelović and A. Zisserman. Three things everyone should know to improve object retrieval. In *IEEE PAMI*, 2012. [4](#)
- [3] K. Chatfield, V. Lempitsky, A. Vedaldi, and A. Zisserman. The devil is in the details: an evaluation of recent feature encoding methods. In *Proc. BMVC*, 2011. [2](#)
- [4] D. Chen, G. Baatz, Köser, S. Tsai, R. Vedantham, T. Pylvanainen, K. Roimela, X. Chen, J. Bach, M. Pollefeys, B. Girod, and R. Grzeszczuk. City-scale landmark identification on mobile devices. In *CVPR*, 2011. [1, 4, 5](#)
- [5] G. Csurka, C. Bray, C. Dance, and L. Fan. Visual categorization with bags of keypoints. In *Workshop on Statistical Learning in Computer Vision, ECCV*, pages 1–22, 2004. [2](#)
- [6] M. Cummins and P. Newman. Highly scalable appearance-only SLAM - FAB-MAP 2.0. In *Proceedings of Robotics: Science and Systems*, Seattle, USA, June 2009. [1](#)
- [7] N. Dalal and B. Triggs. Histogram of oriented gradients for human detection. In *CVPR*, 2005. [7](#)
- [8] C. Doersch, S. Singh, A. Gupta, J. Sivic, and A. A. Efros. What makes paris look like paris? *SIGGRAPH*, 31(4), 2012. [2](#)
- [9] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. *arXiv:1310.1531*, 2013. [7](#)
- [10] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. LIBLINEAR: A library for large linear classification. *J. Machine Learning Research*, 9:1871–1874, 2008. [3, 5](#)
- [11] M. Gharbi, T. Malisiewicz, S. Paris, and F. Durand. A Gaussian approximation of feature space for fast image similarity. Technical report, MIT, 2012. [7](#)

- 756 [12] P. Gronat, G. Obozinski, J. Sivic, and T. Pajdla. Learning and  
757 calibrating per-location classifiers for visual place recogni-  
758 tion. In *CVPR*, 2013. 1, 2, 3, 4, 5 810  
759 [13] B. Hariharan, J. Malik, and D. Ramanan. Discriminative  
760 decorrelation for clustering and classification. In *ECCV*,  
761 2012. 7 811  
762 [14] H. Jégou, M. Douze, and C. Schmid. Product Quantization  
763 for Nearest Neighbor Search. *IEEE PAMI*, 33(1), 2011. 1 812  
764 [15] H. Jégou, F. Perronnin, M. Douze, J. Sánchez, P. Pérez, and  
765 C. Schmid. Aggregating local image descriptors into com-  
766 pact codes. *IEEE PAMI*, 34:1704–1716, 2012. 2, 4, 5 813  
767 [16] B. Klingner, D. Martin, and J. Roseborough. Street view  
768 motion-from-structure-from-motion. In *ICCV*, 2013. 1, 2 814  
769 [17] J. Knopp, J. Sivic, and T. Pajdla. Avoiding confusing features  
770 in place recognition. In *ECCV*, 2010. 1, 5 815  
771 [18] J. Krapac, J. Verbeek, and F. Jurie. Modeling Spatial Lay-  
772 out with Fisher Vectors for Image Categorization. In *ICCV*,  
773 2011. 2 816  
774 [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet  
775 classification with deep convolutional neural networks. In  
776 *NIPS*, 2012. 7 817  
777 [20] Y. Li, N. Snavely, D. Huttenlocher, and P. Fua. Worldwide  
778 pose estimation using 3d point clouds. In *ECCV*, 2012. 2 818  
779 [21] D. Lowe. Distinctive image features from scale-invariant  
780 keypoints. *IJCV*, 60(2):91–110, 2004. 1 819  
781 [22] T. Malisiewicz, A. Gupta, and A. A. Efros. Ensemble of  
782 exemplar-svms for object detection and beyond. In *ICCV*,  
783 2011. 1, 2, 3 820  
784 [23] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and  
785 transferring mid-level image representations using convolu-  
786 tional neural networks. In *CVPR*, 2014. 7 821  
787 [24] F. Perronnin, Y. Liu, J. Snchez, and H. Poirier. Large-scale  
788 image retrieval with compressed fisher vectors. In *CVPR*.  
789 IEEE, 2010. 1 822  
790 [25] T. Sattler, B. Leibe, and L. Kobbelt. Improving image-based  
791 localization by active correspondence search. In *ECCV*,  
792 2012. 2 823  
793 [26] T. Sattler, T. Weyand, B. Leibe, and L. Kobbelt. Image  
794 retrieval for image-based localization revisited. In *Proc.  
795 BMVC*, 2012. 5 824  
796 [27] G. Schindler, M. Brown, and R. Szeliski. City-scale location  
797 recognition. In *CVPR*, 2007. 1 825  
798 [28] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and  
799 Y. LeCun. Overfeat: Integrated recognition, localization and  
800 detection using convolutional networks. *arXiv:1312.6229*,  
801 2013. 7 826  
802 [29] A. Shrivastava, T. Malisiewicz, A. Gupta, and A. A. Efros.  
803 Data-driven visual similarity for cross-domain image match-  
804 ing. In *SIGGRAPH ASIA*, 2011. 2 827  
805 [30] K. Simonyan, O. M. Parkhi, A. Vedaldi, and A. Zisserman.  
806 Fisher Vector Faces in the Wild. In *British Machine Vision  
807 Conference*, 2013. 2 828  
808 [31] S. Singh, A. Gupta, and A. A. Efros. Unsupervised discovery  
809 of mid-level discriminative patches. In *ECCV*, 2012. 2 829  
810 [32] J. Sivic and A. Zisserman. Video Google: A text retrieval  
811 approach to object matching in videos. In *ICCV*, 2003. 2 830