

000
001
002054
055
056

LATEX Fisher vector places: learning compact image descriptors for place recognition

003
004
005
006
007057
058
059
060
061

008 Anonymous CVPR submission

009
010
011062
063

012 Paper ID ****

013
014
015064
065
066

Abstract

016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038067
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

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. Second, 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. Third, 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.

039
040
041
042
043
044
045
046
047
048
049
050
051
052
053094
095
096
097

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 chal-

lenge 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 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.

108	2. Related work	162
109		163
110	Large-scale visual place recognition The visual localization problem is typically treated as large-scale instance-level retrieval [6, 4, 12, 17, 27, 34, 35], where images are represented using local invariant features [21] aggregated into the bag-of-visual-words [5, 32] representation. The image database can be further augmented by 3D point clouds [16], automatically reconstructed by large-scale structure from motion (SfM) [1, 16], which enables accurate prediction of query image camera position [20, 25]. In this work we investigate learning a discriminative representation using the compact Fisher vector descriptors [15]. Fisher vector descriptors have shown excellent place recognition accuracy [34]. In this work we further improve their performance by discriminative learning.	164
111		165
112		166
113		167
114		168
115		169
116		170
117		171
118		172
119		173
120		174
121		175
122		176
123		177
124		178
125	Fisher vector image representations Fisher vector image representations have recently demonstrated excellent performance for a number of visual recognition tasks [3, 15, 18, 30]. They are specially suited for retrieval applications since they are robust to image appearance variations and capture richer image statistics than the simple bag-of-visual-words (BOW) aggregation. However, the raw extracted Fisher vectors are typically high-dimensional, e.g. with 32,768 non-sparse dimensions, which is impractical for large-scale visual recognition and indexing applications. Hence, their dimensionality is often reduced by principal component analysis (PCA) and further quantized for efficient indexing using, e.g. a product quantizer [15]. Other recent work has demonstrated improved performance in a face recognition application by finding discriminative projection using a large amount of training face data [30]. Our work is complementary to these methods as it operates on the projected low-dimensional descriptor and further learns discriminative re-weighting of the descriptor specific to each image in the database using per-exemplar support vector machine [22].	179
126		180
127		181
128		182
129		183
130		184
131		185
132		186
133		187
134		188
135		189
136		190
137		191
138		192
139		193
140		194
141		195
142		196
143		197
144		198
145		199
146		200
147		201
148	Per-exemplar support vector machine The exemplar support vector machine (e-SVM) has been used in a number of visual recognition tasks including category-level recognition [22], cross-domain retrieval [29], scene parsing [33], place recognition [12] or as an initialization for more complex discriminative clustering models [8, 31]. The main idea is to train a linear support vector machine (SVM) classifier from a single positive example and a large number of negatives. The intuition is that the resulting weight vector will give a higher weight to the discriminative dimensions of the positive training data point and will down weight dimensions that are non-discriminative with respect to the negative training data. A key advantage is that each per-exemplar classifier can be trained independently and hence	202
149		203
150		204
151		205
152		206
153		207
154		208
155		209
156		210
157		211
158		212
159		213
160		214
161		215

the learning can be heavily parallelized. The per-exemplar training brings however also an important drawback. As each classifier is trained independently a careful calibration of the resulting classifier scores is required [12, 22].

Contributions The contributions of this work are three-fold. First, 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. Second, we demonstrate that after an appropriate normalization of the new re-weighted descriptor no further calibration is necessary. Third, we apply e-SVM training to compact Fisher vector descriptors for large-scale place recognition resulting in a *discriminative yet compact* representation of each image in the database. Place recognition results are shown on a dataset of 25k images of Pittsburgh and demonstrate the learnt representation consistently improves over the standard Fisher vector descriptors at different target dimensions.

3. Learning compact place descriptors using per-exemplar SVM

Each database image j is represented by its L2-normalized Fisher vector Φ_j . The goal is to learn a set of new L2-normalized Fisher vectors Ψ_j , one per each database image, such that at query time, given the Fisher vector Φ_q of an unknown query image, we retrieve the database image depicting the same location by finding the image j^* with the highest score measured by a dot product

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

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

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

216 where Φ_j^+ is the descriptor of the place j as the positive data
 217 point, Φ^- are Fisher descriptors from negative training data
 218 \mathcal{N}_j and h is the hinge loss, $h(y) = \max(1 - y, 0)$. Note
 219 that the first term in (2) is the regularizer, the second term
 220 is the loss on the positive data weighted by scalar parameter
 221 C_1 and the third term is the loss on the negative data
 222 weighted by scalar parameter C_2 . The objective is convex
 223 and can be minimized with respect to w_j and b_j using stan-
 224 dard software packages such as [10]. A key advantage is
 225 that the per-exemplar classifier for each place can be trained
 226 independently and hence the learning can be heavily paral-
 227 lelized. The downside of the independent training for each
 228 positive example is that the resulting scores have to be cal-
 229 brated with respect to each other on additional data [12, 22].
 230

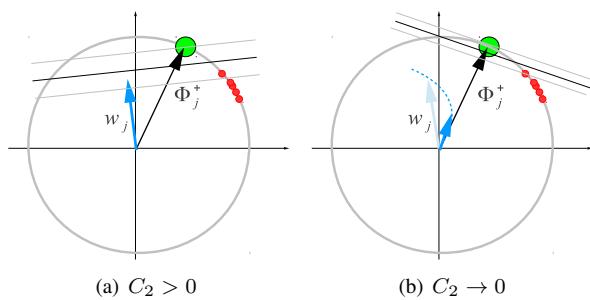
233 **Analysis of per-exemplar SVM objective** In this sec-
 234 tion, we analyze the exemplar SVM objective (2) and show
 235 the learnt and *re-normalized* weight vector w_j can be in-
 236 terpreted as a new descriptor Ψ_j that replaces the original
 237 positive training descriptor Φ_j^+ . In particular, we show first
 238 that when the weight C_2 of the negative data in objective (2)
 239 goes to zero and the learnt Ψ_j is identical to the original
 240 positive training data point Φ_j^+ . Second, when $C_2 > 0$, the
 241 learnt Ψ_j moves away from the positive Φ_j^+ to increase its
 242 separation from the negative data. Details are given next.
 243

246 **Case I:** $C_2 \rightarrow 0$. The goal is to show that when the weight
 247 C_2 of the negative data in objective (2) goes towards zero
 248 the resulting hyperplane vector w_j is parallel with the pos-
 249 itive training descriptor Φ_j^+ . When w_j is normalized to
 250 have unit L2 norm the two vectors are identical. First, let
 251 us decompose w into parallel and orthogonal part with re-
 252 spect to the positive training data point Φ_j^+ (in the following
 253 we omit index j for brevity), i.e. $w = w^\perp + w^\parallel$, where
 254 $(w^\perp)^T \Phi_j^+ = 0$. Next, we observe that when the weight of
 255 the negative data diminishes ($C_2 \rightarrow 0$), any non-zero com-
 256 ponent w^\perp will increase the value of the objective. As a
 257 result, for $C_2 \rightarrow 0$ the objective is minimized by w^\parallel , i.e.
 258 the optimal w is parallel with Φ_j^+ .

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

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

262 Note that the orthogonal part w^\perp does not change the value
 263 of the second term in (3) because $(w^\perp + w^\parallel)^T \Phi_j^+ =$



266 **Figure 1: An illustration of the effect of decreasing pa-
 267 rameter C_2 in the exemplar support vector machine ob-
 268 jective.** The positive exemplar Φ_j^+ is shown in green. The
 269 negative data points are shown in red. All training data is
 270 L2 normalized to lie on a hyper-sphere. (a) For $C_2 > 0$, the
 271 normal w_j of the optimal hyper-plane moves away from the
 272 direction given by the positive example Φ_j^+ in a manner that
 273 reduces the loss on the negative data. (b) As the parameter
 274 C_2 decreases the learnt w_j becomes parallel to the positive
 275 training example Φ_j^+ and its magnitude $\|w_j\|$ goes to 0.

278 $(w^\parallel)^T \Phi_j^+$, and hence (3) reduces to

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

291 In the limit case as $C_2 \rightarrow 0$ any non-zero component w^\perp
 292 will increase the value of the objective (4). This can be
 293 seen by noting that the third term vanishes when $C_2 \rightarrow 0$
 294 and hence the objective is dominated by the first two terms.
 295 Further, the second term in (4) is independent of w^\perp . Fi-
 296 nally, the first term will always increase for any non-zero
 297 value of w^\perp as $\|w^\perp + w^\parallel\|^2 \geq \|w^\parallel\|^2$ for any $w^\perp \neq 0$.

298 As a result, in the limit case when $C_2 \rightarrow 0$ the optimal
 299 w is parallel with Φ_j^+ . Note also, that when C_2 is exactly
 300 equal to zero, $C_2 = 0$, the optimal w vanishes, i.e. the
 301 objective (4) is minimized by trivial solution $\|w\| = 0$ and
 302 $b_j = -1$. The effect of decreasing the parameter C_2 is
 303 illustrated in figure 1.

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

324
 325 **Interpreting normalized w as a new descriptor** The
 326 above analysis demonstrates that as C_2 decreases the
 327 normal of the optimal hyperplane w that separates the positive
 328 exemplar Φ^+ from negative data becomes parallel with Φ^+ ,
 329 as shown in figure 1. As C_2 increases, the normal w of the
 330 optimal hyper-plane moves away from the direction given
 331 by the positive example Φ^+ in a manner that reduces the
 332 loss on the negative data. This suggests that the learnt w
 333 could be interpreted as a modified positive example Φ^+ , re-
 334 weighted to emphasize directions that separate Φ^+ from the
 335 negative data. As discussed above w is not normalized. As
 336 we wish to measure the similarity between descriptors by
 337 (the cosine of) their angle given by equation (1), additional
 338 normalization of the learnt w is necessary. Hence we define
 339 the new descriptor Ψ_j as the normalized hyperplane normal
 340 w_j

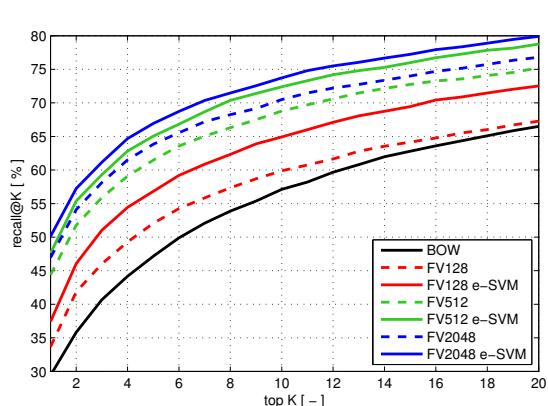
$$\Psi_j = \frac{w_j}{\|w_j\|}. \quad (5)$$

4. Experimental evaluation

344
 345 In this section we first describe the experimental set-up,
 346 then we give implementation details, and finally report re-
 347 sults of the proposed approach on two datasets where we
 348 compare performance with raw Fisher vector matching and
 349 several baselines methods.

4.1. Experimental set-up

350
 351 We perform experiments on a database of Google Street
 352 View images of Pittsburgh downloaded from the Internet.
 353 The data contains panoramas covering roughly an area of
 354 $1.3 \times 1.2 \text{ km}^2$. Similar to [4], for each panorama we
 355 generate 12 overlapping perspective views corresponding
 356 to two different elevation angles to capture both the street-
 357 level scene and the building façades, resulting in a total of
 358 24 perspective views each with 90° FOV and resolution of
 359 960×720 pixels. For evaluation we have used two versions
 360 of this data. The first one was obtained from the authors
 361 of [12] (25k images). In the second version we download
 362 additional images to increase the dataset size to 55k images.
 363 As a query set with known ground truth GPS positions, we
 364 use the 8999 panoramas from the Google Street View re-
 365 search dataset. This dataset covers approximately the same
 366 area, but has been captured at a different time, and depicts
 367 the same places from different viewpoints and under differ-
 368 ent illumination conditions. For each test panorama, we
 369 generate a set of perspective images as described above. Fi-
 370 nally, we randomly select out of all generated perspective
 371 views a subset of 4k images, which is used as a test set to
 372 evaluate the performance of the proposed approach. Since
 373 all the query images have associated GPS location we can
 374 compute their spatial distance from the database images re-
 375 turned by the matching method. We consider a query image
 376 to be correctly localized if the retrieved database image lies
 377



392
 393 Figure 2: **Evaluation on Pittsburgh 25k [12] dataset.** The
 394 fraction of correctly recognized queries (recall@K, y-axis)
 395 vs. the number of top K retrieved database images for
 396 different Fisher vector dimensions. The learnt descriptors
 397 by the proposed method (FV e-SVM) consistently improve
 398 over the raw Fisher vector descriptors across the whole
 399 range of K and all dimensions.

400 within a perimeter of 20m from the location of the query.

4.2. Implementation details

401
 402 We first extract rootSIFT descriptors [2] for each im-
 403 age. Following [15] we project the 128-dimensional SIFT
 404 descriptors to 64 dimensions using PCA. The projection
 405 matrix is learnt on a set of descriptors from 5,000 ran-
 406 domly selected database images. This has also the effect
 407 of decorrelating the SIFT descriptor. The 64-dimensional
 408 SIFT descriptors are then aggregated into Fisher vectors us-
 409 ing a Gaussian mixture model with $N = 256$ components,
 410 which results in a $2 \times 256 \times 64 = 32,768$ -dimensional
 411 descriptor for each image. The Gaussian mixture model
 412 is learnt from descriptors extracted from 5,000 randomly
 413 sampled database images. The high-dimensional Fisher
 414 vector descriptors are then projected down to dimension
 415 $d \in \{128, 512, 2048\}$ using PCA learnt from all avail-
 416 able images in the database. The resulting low dimensional
 417 Fisher vectors are then re-normalized to have unit L2-norm,
 418 which we found to be important in practice.

419
 420 **Learning parameters and training data.** To learn the
 421 exemplar support vector machine for each database image
 422 j , the positive and negative training data are constructed as
 423 follows. The *negative training set* \mathcal{N}_j is obtained by: (i)
 424 finding the set of images with geographical distance greater
 425 than 200m; (ii) sorting the images by decreasing value of
 426 similarity to image j measured by the dot product between
 427 their respective Fisher vectors; (iii) taking the top $N = 500$
 428 ranked images as the negative set. In other words, the neg-
 429 ative training data consists of the hard negative images, i.e.
 430

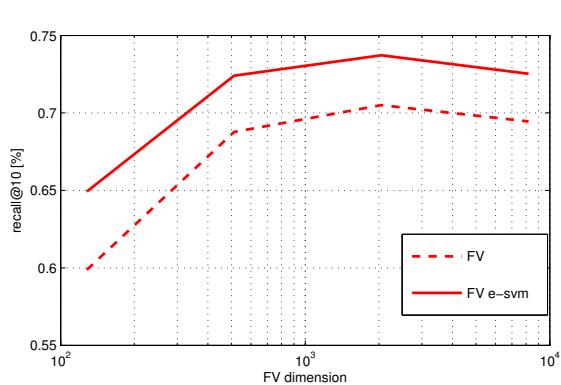
432 those that are similar to image j but are far away from its
 433 geographical position, hence, cannot have the same visual
 434 content. The *positive training set* \mathcal{P}_j consist of the original
 435 Fisher vector Φ_j of the image j . For the SVM training we
 436 use libsvm [10]. We use the same C_1 and C_2 parameters
 437 for all per-exemplar classifiers, but find the optimal value of
 438 the parameters for each dimensionality of the Fisher vector
 439 by a grid search evaluating performance on a held out set.
 440 We observe that for different Fisher vector target dimen-
 441 sions the optimal value of parameter C_1 is quite stable (typi-
 442 cally $C_1 = 1$) while the optimal parameter for C_2 varies
 443 between 10^{-6} to 10^{-1} . To learn the new image representa-
 444 tion for each database image j we: (i) learn SVM from \mathcal{P}_j
 445 and N_j (see above); (ii) $L2$ normalize the learned w_j using
 446 equation (5); and (iii) use this re-normalized vector as the
 447 new image descriptor Ψ_j for image j . At query time we
 448 compute the Fisher vector Φ_q of the query image and mea-
 449 sure its similarity score to the learnt descriptors Ψ_j for each
 450 database image by equation (1).

452 4.3. Results

453 For each database (Pittsburgh 25k and Pittsburgh 55k)
 454 we compare results of our method (FV e-SVM) to two base-
 455 lines: standard bag-of-visual-words baseline (BOW) and
 456 raw Fisher vector matching without learning (FV).
 457

458 We perform experiments on several target Fisher vector
 459 dimensions $d \in \{128, 512, 2048\}$. For each method we
 460 measure performance using the percentage of correctly rec-
 461 gnized queries (Recall) similarly to, e.g., [4, 17, 26]. The
 462 query is correctly localized if at least one of the top K re-
 463 tried database images is within 20 meters from the ground
 464 truth position of the query. Results are shown for different
 465 values of K in table 1. For the Pittsburgh 25k we also show
 466 results in the form of a curve in figure 2. The results clearly
 467 demonstrate the benefits of the learnt descriptors with re-
 468 spect to the standard Fisher vectors for all target dimen-
 469 sions and lengths of shortlist K . The benefits of discrimi-
 470 native learning are specially prominent for low-dimensional
 471 compact descriptors ($d = 128$). The proposed method also
 472 significantly outperforms the bag-of-visual-words baseline.
 473 Figure 4 shows examples of place recognition results.

474
 475 **Comparison to other methods.** On the Pittsburgh 25k
 476 database, we compare performance of our learnt discrimi-
 477 native descriptors to the methods of [12] and [17], who re-
 478 port on the same testing data top $K = 1$ recall of 36.5% and
 479 41.9%, respectively (results taken from [12]). Our method
 480 outperforms [17] already for dimension $d = 128$ (37.8%)
 481 and [12] for dimension $d = 512$ (47.6%). Furthermore,
 482 note that [12, 17] are based on a bag-of-visual-words rep-
 483 resentation, which typically needs to store between 1000-
 484 2000 non-zero visual words per image, which is signifi-
 485 cantly more than our learnt 128 or 512 dimensional descrip-



486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539

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.

508 tor.

512
Memory complexity analysis. Figure 3 compares the
 513 performance of the learnt discriminative descriptor (FV
 514 eSVM) to the raw Fisher vectors (FV) for different target
 515 dimensions. The results demonstrate that for a given level
 516 of accuracy (y-axis) our method learns a more compact
 517 (lower-dimensional) representation (x-axis). For example,
 518 our learnt 128-dimensional descriptor achieves a similar ac-
 519 curacy (around 65%) to the 256-dimensional raw Fisher
 520 descriptor essentially reducing the memory complexity to
 521 50% for the same level of performance. Note that simi-
 522 lar to [15], we observe decrease in performance at high-
 523 dimensions for both the baseline and our method.

5. Conclusions

524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539

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



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.

Method:	25k Pittsburgh					55k Pittsburgh					702
	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	703
FV128	33.6	41.8	52.0	59.8	67.7	10.9	14.1	20.2	26.4	33.2	704
FV128 e-SVM	37.8	46.1	56.9	64.9	72.6	13.5	17.7	25.0	31.8	39.0	705
FV512	44.3	51.7	61.4	68.7	75.2	17.3	21.1	28.4	34.2	40.3	706
FV512 e-SVM	47.6	55.4	65.1	72.4	78.8	19.8	25.1	32.7	38.7	46.0	707
FV2048	46.9	54.1	63.8	70.5	76.8	19.2	23.5	29.9	35.2	41.9	708
FV2048 e-SVM	50.2	57.3	67.0	73.8	78.0	20.8	25.9	33.1	38.7	45.9	709

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.

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. [2, 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. [2](#)
- [7] N. Dalal and B. Triggs. Histogram of oriented gradients for human detection. In *CVPR*, 2005. [5](#)
- [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. [5](#)
- [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](#)
- [12] P. Gronat, G. Obozinski, J. Sivic, and T. Pajdla. Learning and calibrating per-location classifiers for visual place recognition. In *CVPR*, 2013. [1, 2, 3, 4, 5](#)
- [13] B. Hariharan, J. Malik, and D. Ramanan. Discriminative decorrelation for clustering and classification. In *ECCV*, 2012. [7](#)
- [14] H. Jégou, M. Douze, and C. Schmid. Product Quantization for Nearest Neighbor Search. *IEEE PAMI*, 33(1), 2011. [1](#)
- [15] H. Jégou, F. Perronnin, M. Douze, J. Sánchez, P. Pérez, and C. Schmid. Aggregating local image descriptors into compact codes. *IEEE PAMI*, 34:1704–1716, 2012. [2, 4, 5](#)
- [16] B. Klingner, D. Martin, and J. Roseborough. Street view motion-from-structure-from-motion. In *ICCV*, 2013. [1, 2](#)
- [17] J. Knopp, J. Sivic, and T. Pajdla. Avoiding confusing features in place recognition. In *ECCV*, 2010. [2, 5](#)
- [18] J. Krapac, J. Verbeek, and F. Jurie. Modeling Spatial Layout with Fisher Vectors for Image Categorization. In *ICCV*, 2011. [2](#)
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. [5](#)
- [20] Y. Li, N. Snavely, D. Huttenlocher, and P. Fua. Worldwide pose estimation using 3d point clouds. In *ECCV*, 2012. [2](#)
- [21] D. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, 2004. [2](#)
- [22] T. Malisiewicz, A. Gupta, and A. A. Efros. Ensemble of exemplar-svms for object detection and beyond. In *ICCV*, 2011. [1, 2, 3](#)
- [23] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *CVPR*, 2014. [5](#)
- [24] F. Perronnin, Y. Liu, J. Snchez, and H. Poirier. Large-scale image retrieval with compressed fisher vectors. In *CVPR*, IEEE, 2010. [1](#)
- [25] T. Sattler, B. Leibe, and L. Kobbelt. Improving image-based localization by active correspondence search. In *ECCV*, 2012. [2](#)
- [26] T. Sattler, T. Weyand, B. Leibe, and L. Kobbelt. Image retrieval for image-based localization revisited. In *Proc. BMVC*, 2012. [5](#)
- [27] G. Schindler, M. Brown, and R. Szeliski. City-scale location recognition. In *CVPR*, 2007. [2](#)
- [28] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv:1312.6229*, 2013. [5](#)

- 756 [29] A. Shrivastava, T. Malisiewicz, A. Gupta, and A. A. Efros. 810
 757 Data-driven visual similarity for cross-domain image match- 811
 758 ing. In *SIGGRAPH ASIA*, 2011. 2 812
 759 [30] K. Simonyan, O. M. Parkhi, A. Vedaldi, and A. Zisserman. 813
 760 Fisher Vector Faces in the Wild. In *British Machine Vision* 814
 761 Conference, 2013. 2 815
 762 [31] S. Singh, A. Gupta, and A. A. Efros. Unsupervised discovery 816
 763 of mid-level discriminative patches. In *ECCV*, 2012. 2 817
 764 [32] J. Sivic and A. Zisserman. Video Google: A text retrieval 818
 765 approach to object matching in videos. In *ICCV*, 2003. 2 819
 766 [33] J. Tighe and S. Lazebnik. Finding things: Image parsing with 820
 767 regions and per-exemplar detectors. In *CVPR*, 2013. 2 821
 768 [34] A. Torii, J. Sivic, T. Pajdla, and M. Okutomi. Visual place 822
 769 recognition with repetitive structures. In *CVPR*, 2013. 2 823
 770 [35] A. Zamir and M. Shah. Accurate image localization based 824
 771 on google maps street view. In *ECCV*, 2010. 2 825
 772 826
 773 827
 774 828
 775 829
 776 830
 777 831
 778 832
 779 833
 780 834
 781 835
 782 836
 783 837
 784 838
 785 839
 786 840
 787 841
 788 842
 789 843
 790 844
 791 845
 792 846
 793 847
 794 848
 795 849
 796 850
 797 851
 798 852
 799 853
 800 854
 801 855
 802 856
 803 857
 804 858
 805 859
 806 860
 807 861
 808 862
 809 863