Traffic Scene Classification on a Representation Budget

Ivan Sikirić[®], Karla Brkić, Petra Bevandić, Ivan Krešo, Josip Krapac, and Siniša Šegvić, Member, IEEE

Abstract-Visual cues can be used alongside GPS positioning 1 and digital maps to improve understanding of vehicle environ-2 ment in fleet management systems. Such systems are limited both 3 in terms of bandwidth and storage space, so minimizing the size 4 of transmitted and stored visual data is a priority. In this paper, 5 we present efficient strategies for computing very short image representations suitable for classifying various types of traffic 7 scenes in fleet management systems. We anticipate that the set of 8 interesting classes will change over time, so we consider image 9 representations that can be trained without knowing the labels 10 of the target dataset. We empirically evaluate and compare the 11 presented methods on a contributed dataset of 11447 labeled 12 traffic scenes. Our results indicate that excellent classification 13 results can be achieved with very short image representations, 14 and that fine-tuning on the target dataset image data is not 15 mandatory. Image descriptors can be as short as 128 components 16 while still offering good performance, even in presence of adverse 17 weather or illumination conditions. 18

Index Terms—Computer vision, intelligent vehicles, image
 classification.

21

I. INTRODUCTION

N most vision-based systems for scene recognition, size of 22 image representation is of no concern [1]. Very compact 23 image representations are primarily utilized in image retrieval 24 systems [2], but they are not designed to be used for image 25 classification. However, limiting image representation size is 26 critical in systems where images are acquired by thin clients 27 with low processing power and limited bandwidths, as is 28 often the case in intelligent transportation. Many of these 29 systems assume a centralized server which retrieves and stores 30 visual information from a number of thin clients at regular 31 time intervals, as illustrated in Figure 1. In these systems, 32 the number of clients can be very large, while the record 33 keeping time can be very long (several months, or even years). 34 Thus, transferring raw image data from the clients to the server 35

Manuscript received October 16, 2017; revised June 5, 2018 and October 30, 2018; accepted December 14, 2018. This work was supported in part by the Croatian Science Foundation under Grant I-2433-2014 and in part by the European Regional Development Fund (DATACROSS) under Grant KK.01.1.1.01.0009. The Associate Editor for this paper was Q. Ji. (*Corresponding author: Ivan Sikirić.*)

I. Sikirić is with Mireo d.d., 10000 Zagreb, Croatia (e-mail: ivan.sikiric@mireo.hr).

K. Brkić, P. Bevandić, I. Krešo, and S. Šegvić are with the Faculty of Electrical Engineering and Computing, University of Zagreb, 10000 Zagreb, Croatia (e-mail: karla.brkic@fer.hr; petra.bevandic@fer.hr; ivan.kreso@fer.hr; sinisa.segvic@fer.hr).

J. Krapac is with Mobius Labs GmbH, 10997 Berlin, Germany (e-mail: josip@mobius.ml).

Digital Object Identifier 10.1109/TITS.2019.2891995



Fig. 1. The proposed target application framework: a number of thin clients send visual information about their surroundings to a server via limited bandwidth. Figure reproduced from [3].

may be prohibitively expensive both in terms of storage and in terms of data transfer. The latter is especially relevant when thin clients communicate via a paid mobile network, as is often the case in fleet management systems.

Our primary motivation are fleet management systems, where a server tracks the locations of a fleet of vehicles in real time. In most fleet management systems, vehicles are equipped with a GPS sensor and a range of simple supplementary sensors (e.g. ignition sensor, fuel gauge, thermometer). We additionally propose equipping each vehicle with a low cost dashboard camera. The addition of a camera, while not the industry standard in fleet management, would enable the recognition of the type of traffic environment the vehicle is currently in. The traffic environment classification could also be achieved using other sensors, such as radar and lidar [4]. However, in this paper we focus exclusively on visual data for traffic scene classification.

There are several ways to make use of traffic scene clas-53 sification in fleet management systems. One is to aid route 54 reconstruction algorithms, by differentiating between various 55 types of roads. Another is to detect scenarios in which a 56 degradation of GPS accuracy is expected, such as driving 57 through tunnels, under overpasses, toll booths, gas stations, 58 etc. Finally, many companies have their own set of char-59 acteristic locations and traffic scenarios which they would 60 like to detect. These might be specific to their particular 61 business needs, and difficult to guess ahead of time, which is a 62 major obstacle to doing the image classification on the clients 63 themselves. If the classification were to be performed on the 64 client, then client software would have to be updated each 65 time the set of classes is changed. Additionally, it would be 66 impossible to re-classify the historical visual data stored on the 67 server. If the clients simply compute the *image descriptors* and 68 transmit them to the server, then changes in the business logic 69

1524-9050 © 2019 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

89

are localized, cheaper and easier to maintain. Re-examination 70 of historical visual data becomes possible, and the clients are 71 significantly simplified. 72

In this paper we study image classification on a tight repre-73 sentation budget in the context of a thin client - central server 74 scenario. Our focus is on very short image descriptors which 75 stand a fair chance to perform well on classes which have 76 not been seen during training. We consolidate best approaches 77 from our previous study [5] and compare them to recent 78 approaches based on deep convolutional architectures. The 79 performance of all presented approaches is comprehensively 80 evaluated on an extended and improved version of our image 81 classification dataset for fleet management applications. 82

Section II overviews prior research related to this paper. 83 Section III summarizes our approach and categorizes the 84 methods we use for image classification. Section IV details the 85 experimental setup, including the presentation of our dataset, 86 and presents the results. The conclusions and outlook are given 87 in Section V. 88

II. RELATED WORK

Our work is related to three interconnected research areas: 90 (i) image classification in general, (ii) traffic scene classifica-91 tion, and (iii) short image descriptors. 92

A. Image Classification in General 93

Active research in image classification rarely considers 94 descriptor length as a research topic. Rather, the majority of 95 efforts are focused on building reliable methods for classifying 96 a wide range of image categories [6]–[8]. 97

Before the advent of deep learning methods, most success-98 full approaches [9] were based on the seminal bag-of-visual 99 words method [10]. Some extensions improve the performance 100 by augmenting the representation with the spatial layout of 101 visual words [11]. State of the art performance has been 102 obtained with spatial Fisher vectors (SFV) [12] which aggre-103 gate locations of visual words by leveraging Fisher vectors 104 with respect to a spatial generative model. 105

After the success of the AlexNet [13], deep learning meth-106 ods started developing rapidly. Increasing the width and 107 depth of the network further improved classification accuracy, 108 as demonstrated by the architectures VGG-E [14] (19 lay-109 ers) and Inception [15] (22 layers). The 100-layer barrier 110 was surpassed by highway [16] and residual [17] (ResNet) 111 architectures. 112

ResNet reduces the problem of vanishing gradient by 113 introducing shortcut connections between every few stacked 114 layers, thus forming residual building blocks. This architecture 115 substantially improves the convergence of the optimization 116 algorithm. Although this network can be over a hundred layers 117 deep, further research showed that most gradient comes from 118 shallow paths [18]. 119

DenseNet [19] architecture uses dense building blocks in 120 which outputs of every layer are connected to inputs of every 121 subsequent layer in the block. The features are combined 122 via concatenation. The dense connectivity enables feature 123 reuse, simplifies the information flow between layers, requires 124

less parameters per layer than previous approaches, and thus enables training of even greater number of layers.

Some recent approaches have used ImageNet classification 127 as a proxy task for learning high-quality convolutional fea-128 ture extractors [20]. Cimpoi et al. [21] obtain the FV-CNN 129 descriptor by extracting features from deep convolutional 130 layers and aggregating them with a Fisher vector frame-131 work. Their descriptor has achieved state-of-the-art per-132 formance in texture recognition and image classification. 133 Garcia-Gasulla et al. [22] measure the differences of distrib-134 utions of network features in different classes, showing that 135 features of deeper layers are more specialized than features in 136 lower levels. Nanni et al. [23] demonstrate that hand-crafted 137 features and features learned in deep networks extract different 138 information from input images, and that best results may be 139 achieved if both types of features are used, and dimensionality 140 reduced via PCA. 141

The work in this area of research was not limited 142 to supervised learning approaches. Radford et al. [24] 143 extend the ideas of Goodfellow et al. [25] and present 144 an unsupervised learning framework called DCGAN which 145 succeeds to learn image representations from unlabeled 146 data. Arjovsky et al. [26] introduce the WGAN framework, 147 an improvement of DCGAN which results in simplified learn-148 ing process and avoiding the mode collapse problem. 149

B. Traffic Scene Classification

Use of machine vision in vehicles is on the rise [27], 151 so the problem of traffic scene classification is receiving 152 increased attention recently. Oeljeklaus et al. [28] present a 153 deep neural network capable of simultaneous traffic scene 154 recognition and segmentation. The network has less than six 155 million parameters (not counting the FC layer), two separate 156 output paths, and a shared encoder used by both tasks. 157 Di et al. [29], [30] research understanding of traffic scenes 158 from images taken from the same location, but under different 159 weather or illumination conditions. Their approach consists 160 of extracting fine-tuned CNN features and transferring the 161 annotations from the retrieved best matching image based 162 on cross-domain dense correspondences. Hussain et al. [31] 163 demonstrate that adequate vehicle type classification can be 164 done using CNN even if vehicle regions are as small as 165 90×90 pixels. 166

Traffic scene classification is not limited to visual data only. 167 Hsu et al. [32] determine whether a vehicle is driving on 168 roads or viaducts by applying a dynamic Bayesian network 169 on map data, satellite visibility data and position calculated by 170 GPS. Seeger *et al.* [4] demonstrate that four types of traffic 171 scenes can be successfully classified by relying on occupancy 172 grid data. Additionally, they find that using the hand-crafted 173 features with an SVM classifier shows results on par with those achieved by deep learning networks such as VGG16. 175

C. Short Image Descriptors

Most methods for producing short image descriptors are 177 motivated by reducing memory consumption in large-scale 178 image retrieval. They also aim to speed up the retrieval 179

150

125

126

174

process, which is often done via an approximate nearest 180 neighbor search [2]. A common approach is to produce a 181 very short similarity-preserving image hash, such as DSH 182 introduced by Liu et al. [33]. They use supervised information 183 of image pair similarity to train a CNN structure to produce a 184 small number of discrete values on output, ranging in size from 185 12 to 48 bits. Improved results are achieved by approaches 186 that take different image modalities into consideration, such 187 as CSDH [34] and DCH [35]. 188

An alternative to crafting very short descriptors is starting 189 with long descriptors and applying dimensionality reduction 190 methods such as principal component analysis (PCA) [36] or 191 product quantization (PQ) [37], [38]. Both approaches have 192 a potential to detect parts of image representation which are 193 not activated for a particular target dataset [22]. However, this 194 requires training on images from the target dataset, which can 195 be regarded as both an advantage and a disadvantage, as we 196 will discuss in the next section. 197

In our own work [3], [5], we investigate how limiting 198 the size of state-of-the-art image descriptors impacts the 199 performance of traffic scene classification. We propose a 200 short image descriptor that combines compacted spatial Fisher 201 vectors and GIST descriptor in a lossy encoding scheme. 202 Classification performance is retained for the descriptor size 203 as low as 48 bytes per image. The main shortcoming of 204 this method is inadequate performance on some classes of 205 traffic scenes (e.g. dense traffic). This became more appar-206 ent as we expanded our dataset to improve variability of 207 under-represented classes. 208

To the best of our knowledge, no previous work proposes nor evaluates very short image descriptors suitable for classification of traffic scenes, which is the main motivation for this paper.

213

III. OUR APPROACH

We consider short image descriptors (less than 2000 compo-214 nents) which are suitable for multi-label classification. We are 215 motivated by the application scenario where many mobile 216 agents acquire images and transmit image representation over 217 a tight data-channel to a remote server which classifies the 218 representation into an open set of attributes. We explore three 219 different approaches to training of descriptors: (i) unsupervised 220 training on target imagery (addressed in III-A), (ii) supervised 221 transfer learning (addressed in III-B) and (iii) unsupervised 222 transfer learning (addressed in III-C). Descriptors in cate-223 gory (i) are trained on images of traffic scenes, but are unaware 224 of their class labels. Descriptors in category (ii) are trained 225 on an unrelated dataset of labeled images, and then applied 226 to traffic images without any additional training. Descriptors 227 in category (iii) are trained on unrelated dataset of images 228 without class labels, and as such are fully unsupervised. Note 229 that none of the three descriptor categories uses the class 230 labels of traffic scene images during training, which ensures 231 that the set of traffic scene classes remains open. Addition-232 ally, the descriptor categories (ii) and (iii) are not trained 233 on the traffic scene images, which means they are equally 234 likely to perform well on different geographical locations. 235

The descriptors from category (i) might potentially perform better on the geographical location on which they were trained, which would be undesirable. In this paper, all descriptors in categories (ii) and (iii) are trained on the ILSVRC subset of ImageNet [7]. 240

A. Descriptors Based on Unsupervised Training on Target Imagery

Both descriptors in this category are based on the tra-243 ditional bag-of-words framework [10]. The two descriptors 244 use different local features (hand-crafted vs convolutional) 245 and the common aggregation mechanism. The aggregation is 246 performed with spatial Fisher vector embeddings [12], which 247 encode both the appearance and the spatial layout through the 248 use of Fisher vectors. Both approaches require access to target 249 images, since they rely on a visual vocabulary of local features 250 found in the dataset. 251

The first descriptor is a concatenation of GIST scene 252 descriptor [39], [40] and spatial Fisher vectors (SFV) [12] 253 embeddings of SIFT [41] local feature descriptors. The GIST 254 scene descriptor is completely handcrafted and has no training 255 requirements whatsoever. By itself, it performs poorly on 256 our dataset. However, concatenating it to the SFV descriptor 257 improves the performance at very low descriptor lengths. Short 258 image representations can be obtained through the choice of 259 hyper-parameters (such as number of appearance or spatial 260 components), or by producing a large representation first, and 261 subsequently reducing it with PCA. We call this concatenation 262 a SIFT/SFV+GIST descriptor. Of all the descriptors consid-263 ered in this paper, this is the only one that does not make use 264 of ImageNet dataset during training. 265

Next, we replace all hand-crafted features with features 266 obtained with end-to-end learning. More precisely, we drop the 267 GIST part and we replace the SIFT features in spatial Fisher 268 vector framework with responses of a deep convolutional 269 neural network. We use responses of conv5_4 layer of 270 VGG-19 [14] convolutional neural network, trained on the 271 ImageNet dataset. We do not fine-tune the network on our 272 dataset, but the spatial Fisher vector framework does build a 273 visual vocabulary on features found in target images. We refer 274 to the resulting descriptor as VGG/SFV. 275

B. Descriptors Based on Supervised Transfer Learning

Approaches presented in III-A are based on the VGG 277 architecture which aggregates convolutional features with two 278 fully connected layers (fc5, fc6). However, these two layers 279 require over 100 million parameters which makes them very 280 prone to overfitting. Hence, it did not come as a surprise 281 when later research showed that conv5 features aggre-282 gated with Fisher vectors outperform fc6 features in various 283 visual tasks [21]. Recent convolutional architectures [17], [19] 284 replace fully connected layers with simple global average 285 pooling. As a consequence, they succeed to overperform VGG 286 despite using much less parameters. 287

In this category, we consider pooled features extracted by two recent architectures: ResNet [17] and DenseNet [19]. As in the case of VGG, we use public parameterization trained 290

241

242



Fig. 2. Examples of visually degraded images from the FM3a set. (a) Falling rain. (b) Falling snow. (c) Fog. (d) Low sun.

on ImageNet. The resulting features are a good match for our
 task since their dimensionality is rather low while offering a
 high classification potential due to state-of-the-art ImageNet
 performance.

The performance of these features might further increase 295 by performing the aggregation with Fisher vectors since that 296 would present opportunities to leverage spatial layout and/or 297 the knowledge of what is unusual in the target dataset. 298 Nevertheless, we refrain from doing that in order to avoid 299 the descriptors having any knowledge of the target dataset. 300 Consequently, approaches from this category stand a very good 301 chance to generalize well to any geographical region in the 302 world. 303

304 C. Fully Unsupervised Descriptors

327

In this category we consider a convolutional generative 305 adversarial network [24]. More precisely, we work with convo-306 lutional features extracted by the discriminator, whose original 307 task is to detect whether the input image is real or crafted by 308 the generator. In order to perform the classification task, the 309 discriminator network has to learn a sophisticated image rep-310 resentation which we leverage to classify the scene. However, 311 the highest levels of this representation have a much lower 312 semantical quality than in the case of strongly supervised 313 approaches which output a distribution over 1000 classes. 314 We therefore form the image descriptor by concatenating 315 max-pooling evidence from all 6 convolutional layers and 316 regulate its size by adjusting the max-pooling grid. 317

As in section III-B, we use a public parameterization learned 318 on ImageNet [24]. Unlike previously presented supervised 319 models, here the training proceeds without knowing the labels 320 of the training images. In theory, this removes any limita-321 tions due to limited availability of labeled images, and pro-322 vides opportunity to train on billions of web-scraped images. 323 Unfortunately, as we shall see in the experiments, practical 324 performance of fully unsupervised approaches still lags behind 325 the approaches presented in sections III-A and III-B. 326

IV. EXPERIMENTS AND RESULTS

Our experiments evaluate classification performance of the presented image descriptors, with emphasis on very short image representations. The experiments are performed on the novel FM3 dataset¹ which contains labeled traffic scenes of interest for fleet management systems (cf. subsection IV-A). We partition our dataset into fixed train, validation and

¹http://www.zemris.fer.hr/~ssegvic/datasets/unizg-fer-fm3am.tar.gz

test subsets, and employ the common classification setup throughout all experiments (cf. subsection IV-B). We con-335 sider approaches from the three descriptor categories pre-336 sented in Section III and evaluate their performance on 337 our dataset (cf. subsection IV-C-IV-E). We explore the 338 dependency of classification performance on the descriptor 339 length (cf. subsection IV-F). We briefly discuss the results 340 (cf. subsection IV-G), and then proceed to detailed analysis of 341 one of the more successful approaches (cf. subsection IV-H). 342

343

A. The FM3 Dataset

In this work we contribute an improved version of the Fleet 344 Management dataset presented in [5]. The dataset contains 345 11448 images of various traffic scenes in Croatia, as seen 346 from the perspective of a driver. All images are of resolution 347 640×480 , and all were taken during day (from late morning 348 to dusk). The dataset is split into two separate subsets: the 349 main part, dubbed FM3m, which contains 6413 images with 350 good visibility and illumination conditions, and the appendix, 351 dubbed FM3a, which contains 5035 images with various types 352 of visual degradation. Visual degradation in the FM3a set 353 was caused by bad weather conditions (falling rain, falling 354 snow, fog) and/or by very low sun angles (bad illumination 355 due to camera pointing towards sun, and/or low amount of 356 sunlight). Windshield wipers obstruct the scene in 785 of 357 images. Examples of visually degraded images from this set 358 are shown in Figure 2. 359

Most images were captured in the same fashion in which 360 they would be captured in a fleet management system: by an 36 in-car camera, at regular intervals while driving on Croatian 362 roads. Since some interesting types of traffic scenes occur very 363 rarely, we have increased the number of samples for such 364 scenes in the FM3m set by downloading images from Map-365 illary.com. We have hand-picked a batch of images for each 366 rarely occurring class and subsequently manually centered and 367 cropped them to the appropriate point of view, as well as 368 resized them to the resolution of 640×480 pixels. 369

As was discussed in the Introduction, it is not possible to 370 anticipate every type of traffic scene that might be of interest 371 in a fleet management system. The eight classes of FM3 cover 372 a variety of plausible use cases. Class *highway* represents a 373 fast, wide and open road. It must have either two lanes and 374 a shoulder lane, or at least three lanes per direction. This 375 class can be used to help with GPS matching disambiguation, 376 and to verify that fast vehicle traffic is allowed and expected. 377 The class *road* is similar to class *highway*, and represents an 378 open non-highway road outside of settlement. It too can be 379



Fig. 3. Examples of classes from the FM3 dataset. (a) Highway. (b) Road. (c) Tunnel. (d) Exit. (e) Settlement. (f) Overpass. (g) Booth. (h) Traffic.

used to augment GPS matching. The class tunnel represents 380 the entrance to, or inside of a tunnel (but not the very exit). 381 It signals that the loss of GPS signal is expected. The class 382 exit represents the exit of a tunnel, which signals that the 383 GPS signal is expected to return soon. A traffic scene image 384 is labeled as *settlement* if there is visible evidence of vehicle 385 being inside a settlement (e.g. buildings, playgrounds, etc.). 386 This may be interesting for a variety of reasons. For example, 387 GPS precision is often low in settlements. The speed limit is 388 likely to be lower and frequent vehicle stops are expected as 389 well. The class *overpass* signals that the vehicle is directly 390 in front of, or under an overpass. This may help to explain 391 observed loss of GPS precision and to predict that a complete 392 loss of GPS signal is unlikely. Class booth represents a scene 393 directly in front of, or at a toll booth. It also signals a possible 394 loss of GPS precision, and explains the reason of vehicle 395 stopping. This may also be useful for keeping track of travel 396 expenses. Even though most toll booths are present in the 397 map data, the system might easily miss them without the 398 use of visual cues due to temporary loss of GPS precision. 399 Finally the class *traffic* represents the scenes of very dense 400 traffic, or a major occlusion of a scene by a large vehicle. 401 In either case, it explains the low speed and stopping of 402 a vehicle, and as such might be interesting to e.g. deliv-403 ery companies. Example instances of each class are shown 404 in Figure 3. 405

All images were labeled according to these class descrip-406 tions by a single annotator, without the use of visual 407 pre-processing of any kind. Some images have multiple labels 408 409 assigned. We define classes highway, road, tunnel, exit, settle*ment* and *booth* to be mutually exclusive, as they represent a 410 location, but classes overpass and traffic represent an attribute 411 of a scene, and can co-occur with other classes. In some 412 countries it might be logical to allow simultaneous assignment 413 of labels settlement and highway, to detect highways going 414 through large cities, but we had no use of such a scheme in 415 our dataset. The distribution of classes across the FM3m and 416 FM3a sets, and across types of visual degradation is shown 417 in Table I. Note that classes *tunnel*, exit and booth have very 418 few samples in the FM3a set. This is because rain and snow do 419

| TABLE I | |
|-------------------------------------|-----|
| THE DISTRIBUTION OF CLASS LABELS IN | THE |
| FM3m and FM3a Datasets | |

| | FM3m | FM3a | FM3a | FM3a | FM3a | FM3a |
|----------|-------|-------|-----------|--------|--------|-------|
| class | total | total | (low sun) | (rain) | (snow) | (fog) |
| highway | 4646 | 1813 | 1375 | 312 | 12 | 114 |
| road | 390 | 630 | 79 | 401 | 150 | 0 |
| tunnel | 616 | 5 | 2 | 3 | 0 | 0 |
| exit | 73 | 6 | 4 | 1 | 0 | 1 |
| settlem. | 583 | 2572 | 525 | 1707 | 352 | 0 |
| overpass | 152 | 90 | 42 | 32 | 14 | 2 |
| booth | 101 | 5 | 5 | 0 | 0 | 0 |
| traffic | 132 | 351 | 56 | 257 | 38 | 0 |
| total | 6413 | 5035 | 1990 | 2428 | 514 | 115 |

not occur inside tunnels (only on entry and exit), and because 420 class *booth* is a rarely occurring event. 421

B. The Classification Setup

In all experiments classification was performed using an 423 SVM classifier, as it was shown to be a good choice 424 in [3], [21], and [42]. We used LibSVM implementation [43], 425 which was trained in a one-vs-rest fashion. For each of 8 426 binary classifiers, the average precision (AP) is calculated, and 427 mean average precision (mAP) of all 8 classes was used as the 428 measure of classification performance. For each class, 25% of 429 instances were used for training, and another 25% were used 430 for validating the optimal SVM parameters. To counter the 431 bias in the dataset, we weighted the SVM scores based on the 432 inverse proportion of class samples in the training set. Finally, 433 the classifier was trained using these parameters on union 434 of training and validation subsets, and the performance was 435 evaluated on the remaining 50% instances. We experimented 436 with linear and RBF kernel, and in all cases RBF performed 437 better. We do not consider that as a surprise since our image 438 representations are rather small. 439

The FM3a set only contains enough samples for five out of eight classes of interest, and is poorly balanced across types of visual degradation. For that reason, in most experiments we only train and evaluate the methods on images from the FM3m

6

TABLE II Average Precision (%) of SIFT/SFV+GIST Descriptor on FM3m Dataset (SVM With RBF Kernel)

| | descriptor length | | | | | | |
|----------|-------------------|-------|-------|-------|-------|-------|-------|
| class | 1024 | 512 | 256 | 128 | 64 | 32 | 16 |
| highway | 99.73 | 99.81 | 99.83 | 99.85 | 99.81 | 99.7 | 99.34 |
| road | 92.10 | 92.73 | 92.51 | 92.02 | 90.92 | 87.77 | 83.00 |
| tunnel | 99.92 | 99.90 | 99.84 | 99.65 | 99.60 | 99.54 | 99.83 |
| exit | 95.58 | 95.65 | 95.04 | 95.13 | 92.81 | 92.09 | 91.67 |
| settlem. | 96.10 | 96.64 | 96.38 | 96.44 | 95.79 | 94.24 | 92.39 |
| overpass | 96.49 | 95.99 | 96.04 | 94.70 | 92.48 | 86.72 | 87.05 |
| booth | 86.43 | 87.04 | 85.76 | 86.51 | 84.69 | 78.56 | 69.67 |
| traffic | 82.69 | 82.42 | 82.55 | 81.56 | 77.53 | 75.41 | 63.15 |
| mean | 93.63 | 93.77 | 93.49 | 93.23 | 91.70 | 89.25 | 85.76 |

set. We use the FM3a set to test resilience of our approach tovisual degradation of images in subsection IV-H.

446 C. Descriptors With Knowledge of FM3 Image Data

In this category we evaluate descriptors presented in III-A. 447 The SIFT/SFV+GIST descriptor is a concatenation of 448 SIFT/SFV and GIST descriptors. The dimensionality of dense 449 SIFT features d was projected down from d = 128 to 450 d = 80 via PCA. In the SFV framework we used K = 16451 appearance components and C = 1 spatial component, which 452 produces a vector of length K(1 + 2d + 5) = 2656. For 453 GIST we used the implementation provided by Oliva and 454 Torralba [40] without any modifications, which produces a 455 vector of length 512. The final size of the descriptor is thus 456 2656 + 512 = 3168. We have used PCA to further compress 457 the descriptor to lengths 1024, 512 and every other power of 458 two, down to 16. The compressed descriptors were classified 459 with an SVM classifier with the RBF kernel, and the results 460 are shown in Table II. 461

The VGG/SFV descriptor uses convolutional features 462 instead of handcrafted SIFT features. More precisely, it uses 463 the responses of conv5_4 layer of the VGG-19 [14] convo-464 lutional neural network which was trained on the ImageNet 465 dataset. The resolution of our dataset is 640×480 , so the 466 shape of local features is $30 \times 40 \times 512$, i.e. 1200 features of 467 size d = 512. We did not use PCA to reduce dimensionality 468 of these features. As in SIFT/SFV+GIST, we used K = 16469 appearance components and C = 1 spatial component, thus 470 obtaining a vector of length 16480, which was then reduced via 471 PCA to lengths from 1024 to 16. The compressed descriptors 472 were classified with an SVM classifier with the RBF kernel, 473 and the results are shown in Table III. 474

Even though the hand-crafted SIFT/SFV+GIST approach achieves decent results, the deep learning based approach VGG/SFV is clearly better. With only 64 components it outperforms the hand-crafted descriptor of length 1024 (95.9% vs 94.94%).

480 D. Descriptors Trained on ImageNet With Class Labels

In this category we evaluate descriptors presented in III-B. We extract activations from the second-to-last layers obtained with publicly available parameterizations of ResNet-50² [17]

TABLE III Average Precision (%) of VGG/SFV Descriptor on FM3m Dataset (SVM With RBF Kernel)

| | descriptor length | | | | | | |
|----------|-------------------|-------|-------|-------|-------|-------|-------|
| class | 1024 | 512 | 256 | 128 | 64 | 32 | 16 |
| highway | 99.78 | 99.76 | 99.75 | 99.75 | 99.69 | 99.6 | 99.52 |
| road | 91.99 | 91.30 | 91.44 | 90.73 | 89.28 | 87.31 | 85.96 |
| tunnel | 99.93 | 99.94 | 99.95 | 99.95 | 99.96 | 99.92 | 99.91 |
| exit | 96.74 | 96.85 | 97.00 | 96.41 | 96.46 | 95.44 | 94.17 |
| settlem. | 97.49 | 97.46 | 97.48 | 97.38 | 96.87 | 96.31 | 95.28 |
| overpass | 97.55 | 96.95 | 96.11 | 95.96 | 96.49 | 96.03 | 90.22 |
| booth | 99.58 | 99.31 | 98.68 | 99.17 | 98.60 | 97.10 | 89.77 |
| traffic | 85.03 | 84.63 | 85.60 | 84.47 | 82.87 | 77.97 | 74.22 |
| mean | 96.01 | 95.78 | 95.75 | 95.48 | 95.03 | 93.71 | 91.13 |

TABLE IV

Average Precision (%) of ResNet-50 and DenseNet-121 Descriptors on FM3m Dataset (SVM With RBF Kernel)

| | ResNet-50 (length 2048) | | DenseNet-121 (length 1024 | |
|------------|-------------------------|----------|---------------------------|----------|
| class | full FM3m | balanced | full FM3m | balanced |
| highway | 99.84 | 86.73 | 93.68 | 93.60 |
| road | 91.13 | 94.29 | 99.97 | 92.76 |
| tunnel | 99.94 | 99.02 | 97.96 | 99.41 |
| exit | 97.70 | 98.64 | 98.33 | 99.14 |
| settlement | 98.33 | 89.46 | 97.86 | 87.58 |
| overpass | 97.15 | 97.54 | 98.81 | 98.63 |
| booth | 98.75 | 97.91 | 87.85 | 99.57 |
| traffic | 86.75 | 96.34 | 96.80 | 97.74 |
| mean | 96.20 | 94.99 | 96.41 | 96.05 |

and DenseNet-121³ [19] architectures. Therefore, both mod-484 els were trained on labeled ImageNet data, and were not 485 fine-tuned on FM3 dataset. No cropping or resizing of the 486 input images was done. Both feature extractors are fully 487 convolutional and applicable to input images of any resolution. 488 Spatial dimensions of all activation tensors are automatically 489 adjusted to the resolution of the input image. As the second-490 to-last layer in both networks performs global pooling, the two 491 resulting descriptors have length equal to the number of feature 492 maps in that layer, which is 2048 for ResNet-50, and 1024 for 493 DenseNet-121. 494

The obtained classification results are shown in Table IV. 495 We include results obtained on entire FM3m, and a balanced 496 subset containing around 30 training and 30 test images per 497 class. The results indicate that DenseNet-121 descriptor is 498 slightly better, while also being shorter. We try to explain these 499 results as follows. ResNet-50 has three times more parameters 500 than DenseNet-121 (25 vs. 8 million) and performs slightly 501 better on ImageNet. Due to less parameters DenseNet-121 502 is less prone to overfitting than ResNet-50, and therefore 503 achieves better knowledge transfer from ImageNet towards 504 FM3m. To verify this, we have fine-tuned both networks on 505 FM3. Indeed the fine-tuned ResNet-50 achieved better perfor-506 mance than fine-tuned DenseNet-121 (98.0% vs 97.1% AP). 507

E. Descriptors Trained on ImageNet Without Class Labels

In this category we evaluate the descriptor presented in III-C ⁵⁰⁹ by applying the public DCGAN parameterization to input ⁵¹⁰ images downsampled to 32×32 pixels [24]. We first extract ⁵¹¹



Fig. 4. Average precision (%) of selected descriptors on FM3m dataset with respect to representation budget. Particular representations are obtained via PCA as necessary. Mean AP values using linear SVM classifier are shown on the left. Mean AP values using SVM with RBF kernel are shown in the middle. AP values for the class *traffic* using SVM with RBF kernel are shown on the right. Best viewed in color. (a) Linear SVM, all classes. (b) SVM with RBF kernel, all classes. (c) SVM with RBF kernel, class traffic.

TABLE V

AVERAGE PRECISION (%) OF DESCRIPTORS BASED ON DCGAN DISCRIMINATOR ON FM3M DATASET (SVM CLASSIFIER)

| trained on | ImageNet | ImageNet | ImageNet | untrained |
|---------------|--------------|--------------|--------------|--------------------|
| kernel | linear | RBF | linear | linear |
| max-pool grid | 4×4 | 1×1 | 1×1 | 1×1 |
| length | 28672 | 1792 | 1792 | 1792 |
| highway | 99.71 | 99.86 | 99.37 | 99.47±0.11 |
| road | 93.63 | 92.10 | 81.32 | 84.79 ± 1.72 |
| tunnel | 99.07 | 99.33 | 99.48 | $99.20 {\pm} 0.11$ |
| exit | 98.53 | 93.67 | 98.14 | $97.73 {\pm} 0.61$ |
| settlement | 96.74 | 95.43 | 92.91 | 89.31±1.33 |
| overpass | 81.67 | 83.85 | 83.08 | 78.44 ± 3.56 |
| booth | 88.59 | 86.53 | 86.24 | 71.42 ± 3.23 |
| traffic | 78.53 | 66.47 | 65.28 | 62.92 ± 3.29 |
| mean | 92.06 | 89.66 | 88.23 | 85.41 ± 0.81 |

discriminator features by 4×4 max-pooling on each layer, and 512 concatenating the results into a descriptor of length 28672. 513 In the subsequent experiments, we shorten the descriptor by 514 reducing the size of the max-pooling grid to 1×1 . This results 515 in descriptors of length 1792. We also include the result of 516 untrained discriminator initialized with random weights. This 517 was done to test how much knowledge transfer happens from 518 ImageNet to FM3, and how much descriptiveness is due to 519 network structure itself. We repeat this with 1000 different 520 random initializations, and report the mean and the standard 521 deviation. All these results are shown in Table V. 522

Much of the descriptiveness seems to come from the convolutional structure itself, since random-initialized discriminator achieves mAP of 85%. Some performance is lost when 4×4 max-pooling is replaced with 1×1 max-pooling, but 16-fold reduction of descriptor size makes this drop in performance an acceptable trade-off.

529 F. Impact of Descriptor Length

To measure the impact of descriptor length reduction on 530 classification performance, we use PCA to compress the size 531 of select descriptors to a series of different lengths. Starting 532 from 1024, we consider all lengths that are a power of two, 533 down to 16. The SIFT/SFV+GIST and VGG/SFV descriptors 534 already use PCA as a part of their framework, and their full 535 results are shown in Tables II and III. We now apply the PCA 536 to ResNet-50, DenseNet-121 and DCGAN 4×4 descriptors as 537 well. Since PCA is used on vectors calculated on FM3 dataset, 538 all compressed descriptors fall into the category of methods 539



Fig. 5. ROC curves for the SVM classifier with RBF kernel trained on the DenseNet-121 descriptor of length 128. Measured on the FM3m dataset. Best viewed in color.

with knowledge of FM3 image data. Note that full length of 540 DenseNet-121 descriptor is 1024, so PCA was only used to 541 produce its representations of length 512 and lower. 542

We classify the descriptors both with the linear SVM classi-543 fier and the SVM with RBF kernel. The summary of the results 544 is best seen in graph form, in Figure 4 (a, b), which shows the 545 mean average precision of all classes for each of the methods, 546 with respect to the length of the representation, separated for 547 linear and RBF kernels. The RBF kernel is strictly better here, 548 but it is not as noticeable for large representations as it is for 549 very small ones. We conclude that RBF is very advantageous 550 when the representations are very small, although the linear 551 SVM may prove the only applicable solution for large training 552 datasets. 553

One thing to note is that reducing the size of the 554 representation does not change the relative order of methods. 555 The DenseNet-121 is a clear winner, followed by ResNet-556 50 and VGG/SFV, followed by SIFT/SFV+GIST (which relies 557 on hand-crafted features) and finally the DCGAN 4×4 (which 558 is trained in an unsupervised manner). Another thing to note 559 is that reducing the size from 1024 down to 128 barely shows 560 any changes in the mAP for most methods if RBF kernel is 561 used. The drop is noticeable for lengths 64 and 32, and very 562 noticeable for length 16. 563

We are not only interested in mean average precision. ⁵⁶⁴ Ideally, every class in the dataset should be classified at ⁵⁶⁵ acceptable levels of error. All descriptors show the worst ⁵⁶⁶ performance on the *traffic* class, so the results for that class are ⁵⁶⁷ shown in Figure 4 (c). Note that AP axis starts from 50% and ⁵⁶⁸ that a considerable drop with respect to the mean performance ⁵⁶⁹ can be observed for all representation budgets. ⁵⁷⁰



Fig. 6. Examples of mispredictions for the SVM classifier with RBF kernel on DenseNet-121 descriptors of length 128. Examples (a) to (f) are from the FM3m testing set, while examples (g) and (h) are from the FM3a set. In all examples the SVM was trained on the FM3m training set. (a) Road as highway. (b) Road as highway. (c) Settlement as highway. (d) Road as settlement. (e) Traffic false positive. (f) Traffic false negative. (g) Settlement as road (rain). (h) Highway as road (fog).

571 G. Discussion

Some classes in the FM3 dataset are clearly much easier 572 to classify than others. In fact, the relative order of classes 573 with respect to their achieved average precision is very similar 574 across all presented methods. The best results are obtained on 575 the highway class which is the most represented class in the 576 dataset. The *tunnel* class is easy to classify, likely because 577 tunnel scenes are almost always dominated by black or orange 578 colors. The *road* class is most often confused with *highway* 579 class. It is also less represented, and has greater variability of 580 visual appearance. By far the hardest class, across all methods, 581 seems to be the traffic class, which includes both the scenes of 582 very dense traffic (which can happen in any kind of location), 583 and also images with major occlusions of scene by other 584 vehicles. One other thing to note is that the class booth is 585 only hard to SIFT/SFV+GIST and DCGAN descriptors, while 586 other methods have excellent performance on this class. 587

We can conclude that the best results are achieved by 588 supervised deep convolutional models, even the ones that 589 have no knowledge of the FM3 dataset. In fact, the best 590 performing method is DenseNet-121, with descriptor length 591 of 1024 and no knowledge of the target dataset. Fine-tuning 592 these approaches to the FM3 dataset would further improve 593 those results, but we avoid that in order to prevent overfitting. 594 DCGAN is at a severe disadvantage compared to supervised 595 approaches. However, it could be improved by training on 596 images acquired world-wide. Such training would not depend 597 on the class labels, which is very desirable in our proposed 598

600 H. Analysis of the DenseNet-121 Descriptor

framework.

599

We now analyze one of our best performing short descriptors: the DenseNet-121 descriptor of length 128. It achieves classification performance almost equal to that of full-length descriptor while being much shorter. We consider it as a strong candidate for our proposed fleet management setup. Its confusion matrix is shown in Table VI. Note that we only include the six classes that are defined as mutually exclusive.

TABLE VI Confusion Matrix for DenseNet-121 Descriptor of Length 128, SVM Classifier With RBF Kernel

| | | | predicted | class | | |
|--------------|---------|------|-----------|-------|--------|-------|
| actual class | highway | road | tunnel | exit | settl. | booth |
| highway | 2298 | 2 | 0 | 0 | 6 | 0 |
| road | 28 | 153 | 0 | 0 | 5 | 0 |
| tunnel | 0 | 0 | 302 | 3 | 0 | 0 |
| exit | 2 | 0 | 1 | 37 | 0 | 0 |
| settlement | 16 | 4 | 0 | 0 | 276 | 0 |
| booth | 1 | 0 | 0 | 0 | 0 | 44 |

The ROC curves for all classes are shown in Figure 5. Some examples of mispredicted classes are shown in Figure 6.

I. Resilience of DenseNet-121 to Visual Degradation

We now use the FM3a set of images to verify the resilience 611 of DenseNet-121 descriptor of length 128 to various types of 612 visual degradation. First we evaluate how the SVM classifier 613 trained on the FM3m training set performs on the FM3a set 614 of degraded images. Next, we add some of images from the 615 FM3a to the training set, and test the performance on the rest 616 of images from FM3a. We do this in two steps, first adding 617 only 10% of images from FM3a, then adding 25% of images 618 to the training set. The average precision results are shown 619 in Table VII. Note that we only list the results for the five 620 classes that are adequately represented in FM3a, excluding 621 classes tunnel, exit and booth. Even without any training on the 622 degraded images, the system was able to classify the degraded 623 images surprisingly well, with the exception of class road, 624 which was most commonly confused with classes settlement 625 (93 times) and highway (82 times). By adding some degraded 626 images to the training set, the quality of classification rises 627 quickly up to 94.46%, which is very close to mAP value for 628 the same classes on non-degraded images: 95.22%. The ROC 629 curves for the five evaluated classes are shown in Figure 7. 630

V. CONCLUSION AND OUTLOOK

In this work we have considered image classification of traffic scenes under requirements specific to fleet management applications, namely i) the bandwidth should be used sparingly 634

610



ROC curves for the SVM classifier with RBF kernel trained on the DenseNet-121 descriptor of length 128. Training data contains 50% of FM3m Fig. 7. samples, and incrementally increasing amounts of FM3a samples. Tested on the rest of FM3a images. Best viewed in color. (a) Training set contains no FM3a samples. (b) Training set contains 10% of FM3a samples. (c) Training set contains 25% of FM3a samples.

TABLE VII AVERAGE PRECISION (%) OF DENSENET-121 DESCRIPTORS OF LENGTH 128 ON FM3A DATASET (SVM WITH RBF KERNEL)

| | not trained on | trained on 10% | trained on 25% |
|------------|----------------|----------------|----------------|
| class | FM3a images | of FM3a images | of FM3a images |
| highway | 85.25 | 98.42 | 99.55 |
| road | 18.87 | 83.52 | 90.30 |
| settlement | 95.73 | 98.04 | 98.85 |
| overpass | 72.08 | 76.97 | 90.07 |
| traffic | 71.70 | 86.96 | 93.53 |
| mean | 68.73 | 88.78 | 94.46 |

to avoid excessive costs, and ii) the number of image classes 635 has to be open. To satisfy these requirements, we studied 636 image descriptors with a restricted representation budget and 637 with no knowledge of the target dataset labels. We achieved 638 resistance to overfitting by considering the following descrip-639 tor training requirements: i) no knowledge of the target dataset 640 labels, ii) no knowledge of the target dataset image data, and 641 iii) no use of supervised training whatsoever. 642

Our experiments have empirically compared a range of clas-643 sification approaches. We have considered handcrafted image 644 descriptors (GIST, SIFT), non-linear embeddings with respect 645 to placement and appearance distribution of image patches 646 (spatial Fisher vectors), and state-of-the-art convolutional rep-647 resentations (VGG, DenseNet, ResNet and DCGAN) trained 648 on ImageNet. Additionally, we have used PCA to reduce 649 image representations down to as low as 16 components, and 650 investigated the drop of classification performance with respect 651 652 to representation length.

Best performance is achieved by deep convolutional models 653 trained in a supervised manner, followed by handcrafted 654 models, and finally by a completely unsupervised descriptor 655 based on DCGAN. We have shown that excellent performance 656 can be achieved even with methods that have no knowledge 657 of our target dataset. With an adequate classifier, image repre-658 sentations can be extremely reduced via PCA (down to as few 659 as 128 components), while sacrificing negligible classification 660 performance. In some cases, even representations with as few 661 as 32 components provide useful results. Finally, we have 662 shown that the best performing method, the DenseNet descrip-663 tor, performs well even on images with visual degradation 664 caused by bad weather and low sun angles, provided some of 665 the degraded images are added to the training data. 666

Future work will include further compression of convo-667 lutional descriptors by fine-tuning on ImageNet, advanced 668

generative adversarial models, and exploring the generalization 669 across datasets taken in different countries.

REFERENCES

- [1] J. Sanchez, F. Perronnin, T. Mensink, and J. Verbeek, "Image classification with the fisher vector: Theory and practice," Int. J. Comput. Vis., vol. 105, no. 3, pp. 222-245, 2013.
- [2] J. Wang, T. Zhang, J. Song, N. Sebe, and H. T. Shen, "A survey on learning to hash," IEEE Trans. Pattern Anal. Mach. Intell., vol. 40, no. 4, pp. 769-790, Apr. 2018.
- [3] I. Sikirić, K. Brkić, J. Krapac, and S. Šegvić, "Robust traffic scene recognition with a limited descriptor length," in Proc. CVPR Workshop VPRCE, 2015, pp. 1-9.
- C. Seeger, A. Müller, L. Schwarz, and M. Manz, "Towards road type [4] classification with occupancy grids," in Proc. IEEE Intell. Vehicles Symp. Workshop, Jun. 2016, pp. 1-4.
- [5] I. Sikirić, K. Brkić, J. Krapac, and S. Šegvić, "Image representations on a budget: Traffic scene classification in a restricted bandwidth scenario," in Proc. IEEE Intell. Vehicles Symp. Workshop, Jun. 2014, pp. 845-852.
- [6] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The PASCAL visual object classes (VOC) challenge," Int. J. Comput. Vis., vol. 88, no. 2, pp. 303-338, 2010.
- [7] O. Russakovsky et al., "ImageNet large scale visual recognition challenge," Int. J. Comput. Vis., vol. 115, no. 3, pp. 211-252, 2015.
- [8] A. Krizhevsky, "Learning multiple layers of features from tiny images," Canadian Inst. Adv. Res., Toronto, ON, Canada, Tech. Rep., 2009.
- K. Chatfield, V. Lempitsky, A. Vedaldi, and A. Zisserman, "The devil [9] is in the details: An evaluation of recent feature encoding methods," in Proc. BMVC, 2011, pp. 1-8.
- [10] G. Csurka, C. R. Dance, L. Fan, J. Willamowski, and C. Bray, "Visual categorization with bags of keypoints," in Proc. Workshop Stat. Learn. Comput. Vis. (ECCV), 2004, pp. 1-2.
- [11] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in Proc. CVPR, Jun. 2006, pp. 2169-2178.
- [12] J. Krapac, J. J. Verbeek, and F. Jurie, "Modeling spatial layout with Fisher vectors for image categorization," in Proc. ICCV, Nov. 2011, pp. 1487-1494.
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in Proc. Adv. Neural Inf. Process. Syst., vol. 25, 2012, pp. 1097-1105.
- [14] K. Simonyan and A. Zisserman. (2014). "Very deep convolutional networks for large-scale image recognition." [Online]. Available: https:// arxiv.org/abs/1409.1556
- C. Szegedy et al., "Going deeper with convolutions," in Proc. CVPR, [15] Jun. 2015, pp. 1-9.
- [16] R. K. Srivastava, K. Greff, and J. Schmidhuber, "Training very deep networks," in Proc. 28th Int. Conf. Neural Inf. Process. Syst. (NIPS), 2015, pp. 2377-2385.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proc. CVPR, Jun. 2016, pp. 770-778.
- [18] A. Veit, M. J. Wilber, and S. J. Belongie, "Residual networks behave like ensembles of relatively shallow networks," in Proc. 29th Annu. Conf. Neural Inf. Process. Syst., 2016, pp. 550-558.
- [19] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in Proc. CVPR, Jul. 2017, pp. 1-3.
- M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Learning and transferring [20] mid-level image representations using convolutional neural networks. in Proc. CVPR, Jun. 2014, pp. 1717-1724.

- M. Cimpoi, S. Maji, I. Kokkinos, and A. Vedaldi, "Deep filter banks for texture recognition, description, and segmentation," *Int. J. Comput. Vis.*, vol. 118, no. 1, pp. 65–94, 2016.
- [22] D. Garcia-Gasulla *et al.* (2017). "On the behavior of convolutional nets for feature extraction." [Online]. Available: https://arxiv.org/abs/1703.01127
- [23] L. Nanni, S. Ghidoni, and S. Brahnam, "Handcrafted vs. non-handcrafted features for computer vision classification," *Pattern Recognit.*, vol. 71, pp. 158–172, Nov. 2017.
- [24] Å. Radford, L. Metz, and S. Chintala. (2015). "Unsupervised representation learning with deep convolutional generative adversarial networks."
 [Online]. Available: https://arxiv.org/abs/1511.06434
- [25] I. Goodfellow *et al.*, "Generative adversarial nets," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 27, 2014, pp. 2672–2680.
- [26] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *Proc. Mach. Learn. Res. (ICML)*, 2017, pp. 214–233.
- [27] B. Ranft and C. Stiller, "The role of machine vision for intelligent vehicles," *IEEE Trans. Intell. Vehicles*, vol. 1, no. 1, pp. 8–19, Mar. 2016.
- [28] M. Oeljeklaus, F. Hoffmann, and T. Bertram, "A combined recognition and segmentation model for urban traffic scene understanding," in *Proc. ITSC*, Oct. 2017, pp. 1–6.
- [29] S. Di, H. Zhang, X. Mei, D. Prokhorov, and H. Ling, "A benchmark for cross-weather traffic scene understanding," in *Proc. ITSC*, Nov. 2016, pp. 2150–2156.
- [30] S. Di, H. Zhang, C.-G. Li, X. Mei, D. Prokhorov, and H. Ling,
 "Cross-domain traffic scene understanding: A dense correspondencebased transfer learning approach," *IEEE Trans. Intell. Transp. Syst.*,
 vol. 19, no. 3, pp. 745–757, Mar. 2018.
- K. F. Hussain, M. Afifi, and G. Moussa, "A comprehensive study of the effect of spatial resolution and color of digital images on vehicle classification," *IEEE Trans. Intell. Transp. Syst.*, to be published.
- [32] L.-T. Hsu, Y. Gu, and S. Kamijo, "Intelligent viaduct recognition and driving altitude determination using gps data," *IEEE Trans. Intell. Vehicles*, vol. 2, no. 3, pp. 175–184, Sep. 2017.
 [33] H. Liu, R. Wang, S. Shan, and X. Chen, "Deep supervised hashing for
- [33] H. Liu, R. Wang, S. Shan, and X. Chen, "Deep supervised hashing for fast image retrieval," in *Proc. CVPR*, Jun. 2016, pp. 2064–2072.
- [34] L. Liu, Z. Lin, L. Shao, F. Shen, G. Ding, and J. Han, "Sequential discrete hashing for scalable cross-modality similarity retrieval," *IEEE Trans. Image Process.*, vol. 26, no. 1, pp. 107–118, Jan. 2017.
- [35] X. Xu, F. Shen, Y. Yang, H. T. Shen, and X. Li, "Learning discriminative binary codes for large-scale cross-modal retrieval," *IEEE Trans. Image Process.*, vol. 26, no. 5, pp. 2494–2507, May 2017.
- [36] I. Jolliffe, *Principal Component Analysis*. New York, NY, USA:
 Springer-Verlag, 1986.
- [37] H. Jegou, M. Douze, and C. Schmid, "Product quantization for nearest neighbor search," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 1, pp. 117–128, Jan. 2011.
- T. Ge, K. He, Q. Ke, and J. Sun, "Optimized product quantization,"
 IEEE Trans. Pattern Anal. Mach. Intell., vol. 36, no. 4, pp. 744–755,
 Apr. 2014.
- [39] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *Int. J. Comput. Vis.*, vol. 42, no. 3, pp. 145–175, 2001.
- [40] A. Oliva and A. B. Torralba, "Scene-centered description from spatial envelope properties," in *Proc. BMCV*, 2002, pp. 263–272.
- [41] D. G. Lowe, "Distinctive image features from scale-invariant keypoints,"
 Int. J. Comput. Vis., vol. 60, no. 2, pp. 91–110, 2004.
- [42] A. Mahmood, M. Bennamoun, S. An, and F. Sohel. (2016). "ResFeats: Residual network based features for image classification." [Online].
 Available: https://arxiv.org/abs/1611.06656
- 787 [43] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," ACM Trans. Intell. Syst. Technol., vol. 2, no. 3, p. 27, 2011.



Karla Brkić received the M.S. and Ph.D. degrees 797 in computer science from the University of Zagreb, 798 Croatia. She has been with the Faculty of Electrical 799 Engineering and Computing, University of Zagreb, 800 since 2008, where he is currently a Post-Doctoral 801 Researcher. Her scientific interests include traffic 802 sign detection and recognition, spatio-temporal data 803 representations, action recognition, and machine 804 learning. 805



Petra Bevandić received the B.Sc. and M.Sc. 806 degrees from the Faculty of Electrical Engineering 807 and Computing, Zagreb, Croatia, in 2012 and 2014, 808 respectively. She is currently a Teaching Assistant 809 with the Faculty of Electrical Engineering and Com-810 puting, Department of Electronics, Microelectronics, 811 Computer and Intelligent Systems, University of 812 Zagreb. 813



Ivan Krešo received the B.Sc. and M.Sc. degrees 814 in computer science from the Faculty of Electrical 815 Engineering and Computing, Zagreb, in 2011 and 816 2013, respectively. He is currently pursuing the 817 Ph.D. degree with the Department of Electronics, 818 Microelectronics, Computer and Intelligent Systems, 819 University of Zagreb. He is also with the Faculty 820 of Electrical Engineering and Computing, Univer-821 sity of Zagreb, as a Research Engineer with the 822 Department of Electronics, Microelectronics, Com-823 puter and Intelligent Systems. His research interests 824

825

are in the areas of computer vision and machine learning.



Josip Krapac received the B.S. and M.S. degrees 826 in computer science from the University of Zagreb. 827 Croatia, and the Ph.D. degree from the Université de 828 Caen Basse-Normandie, France. He spent a year and 829 a half as a Post-Doctoral Researcher with INRIA, 830 Rennes, France, from 2011 to 2013, and three and 831 a half years, as a Post-Doctoral Researcher with 832 the University of Zagreb. He is currently a Senior 833 Computer Vision Scientist with Mobius Labs GmbH. 834 His research interests include image representations 835 for object classification, detection and segmentation, 836 837

and learning models from minimal amount of supervision.



Ivan Sikirić received the M.S. degree in computer science from the Faculty of Electrical Engineering and Computing, University of Zagreb, Croatia, where he is currently pursuing the Ph.D. degree. Since 2008, he has been with Mireo d.d. His professional and scientific interests include computer vision, fleet management, and navigation and intelligent transportation systems.



Siniša Šegvić received the B.S., M.S., and Ph.D. 838 degrees in electrical engineering and computer 839 science from the University of Zagreb, Croatia. 840 He spent one year as a Post-Doctoral Researcher 841 with IRISA/INRIA, Rennes, France, and also with 842 TU Graz, Austria. He is currently a Full Profes-843 sor with the Faculty of Electrical Engineering and 844 Computing, University of Zagreb. His research and 845 professional interests focus on deep convolutional 846 architectures for classification, semantic segmenta-847 tion, and stereo reconstruction. 848