



HAL
open science

Evaluation of Off-The-Shelf CNNs for the Representation of Natural Scenes with Large Seasonal Variations

Amandine Gout, Yann Lifchitz, Titouan Cottencin, Quentin Groshens,
Jérémy Fix, Cédric Pradalier

► **To cite this version:**

Amandine Gout, Yann Lifchitz, Titouan Cottencin, Quentin Groshens, Jérémy Fix, et al.. Evaluation of Off-The-Shelf CNNs for the Representation of Natural Scenes with Large Seasonal Variations. [Research Report] UMI 2958 GeorgiaTech-CNRS; CentraleSupélec UMI GT-CNRS 2958 Université Paris-Saclay. 2017. hal-01448091v1

HAL Id: hal-01448091

<https://hal.science/hal-01448091v1>

Submitted on 27 Jan 2017 (v1), last revised 9 Feb 2017 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Evaluation of Off-The-Shelf CNNs for the Representation of Natural Scenes with Large Seasonal Variations

Amandine Gout^{1,3}, Yann Lifchitz^{1,3}, Titouan Cottencin^{1,3}, Quentin Groshens^{1,3},
Shane Griffith^{2,3,4}, J  r  my Fix^{1,4}, C  dric Pradalier^{2,3,4}

Abstract—This paper focuses on the evaluation of deep convolutional neural networks for the analysis of images of natural scenes subjected to large seasonal variation as well as significant changes of lighting conditions. The context is the development of tools for long-term natural environment monitoring with an autonomous mobile robot.

We report various experiments conducted on a large dataset consisting of a weekly survey of the shore of a small lake over two years using an autonomous surface vessel. This dataset is used first in a place recognition task framed as a classification problem, then in a pose regression task and finally the internal features learned by the network are evaluated for their representation power.

All our results are based on the Caffe library and default network structures where possible.

I. INTRODUCTION



Fig. 1. Examples of variation in appearance of a section of the lake shore from winter to summer. The significant variation in the vegetation and lighting conditions makes place recognition particularly challenging.

Long-term natural environment monitoring using visual inspection is the process of collecting images of an outdoor landscape over a long duration with respect to the natural dynamics of this environment. In our case, we are specifically considering the weekly observation of a lake shore over multiple years using an autonomous boat programmed to follow the shore at a constant distance while recording images. As such, our images depict scenes which are combinations of close-up trees, far-away trees, bushes, lawns, water and sky, with a high level of similarities. Because we work in a natural setting, this environment is subjected to seasonal changes (trees blooming, leaves falling, ...), structural changes (cut branch, fallen trees, mowed lawns) and weather variations impacting the lighting condition, spectrum and incidence. Figure 1 gives an example of a relatively easy group of images in our dataset.

To assess the difficulty of interpreting these images, in a previous work, we evaluated the time required by a human

to decide if two images correspond to the same place albeit at different time of the year. For some image pairs, our test subjects took more than 30 seconds to confirm their answer.

One of the challenges of natural environment monitoring is to be able to compare the appearance of vastly varying outdoor settings for detecting and classifying changes (e.g., structural damage).

In this context, this paper focuses on the evaluation of deep convolutional neural networks (CNN, [1]) applied to images of natural environments subjected to large seasonal variations. The task we set ourselves has three stages. For the first, we evaluate the ability of a CNN to recognize a place independently of the seasonal and lighting changes. This problem is framed as a classification task where each class is a location around our lake shore and a standard network structure can be used. In the second stage, we consider a CNN trained to predict the pose (or view point) of the camera (2D position and heading) using an adapted network structure suitable for a regression task. Finally, in the third stage, we evaluate the quality of the internal representation learned by the convolutional layers of the CNN to describe a place independently of its seasonal appearance as well as the generalizability of the resulting features. This third stage is important for the potential of this internal representation to be used to detect changes with respect to a season-invariant representation of a place.

This study serves not to develop new network architectures, but rather to evaluate standard off-the-shelf tools in the context of a recurring question in robotics. For this reason, all the results presented in this paper have been trained using the Caffe library [2] with a default network structure where possible.

In summary, this paper contribution is two-fold. First, it introduces a relatively unique long-term natural environment monitoring dataset. Second, this paper is a benchmark on the performance of off-the-shelf CNNs for the very particular task of processing images of natural environments subjected to large seasonal changes and natural lighting conditions variation.

II. DATASET AND EXPERIMENTAL SETUP

We have been creating a dataset of a natural environment since August 2013, which is the basis for this study. Since then, every 8 to 15 days, we operate an autonomous surface vessel (Kingfisher from Clearpath Robotics, see fig. 2) around a small lake next to our campus. This lake is 400m long by 200m wide, with a small island and a total perimeter

¹ Centrale Supelec, Metz, France

² College of Computing, Georgia Institute of Technology, Atlanta, USA

³ GeorgiaTech Lorraine, Metz, France

⁴ CNRS UMI 2958 GT-CNRS, Metz, France

of 1km. Its shores are covered with trees, from small bushes to tall full-grown trees, some at the water line, others further away, grass areas are mixed with the shrubberies and a small scenic trail runs around the lake. Some of the places (as in Fig. 1) have office buildings in the background.



Fig. 2. The Kingfisher on a very smooth lake. The pan-tilt camera is housed in the white dome at the back of the boat, just behind the laser scanner used for navigation.

Every survey we collect contains images acquired by a pan-tilt surveillance camera (704×480 pixels) at 10Hz with a conservative JPEG compression. The boat runs autonomously at a constant distance of 10m to the shore (lattice-type local planner) and at a bit less than 1km/h. This means that a survey is a collection of close to 40'000 images, acquired with the camera pointing to the port or starboard side of the boat (i.e. $\pm 90^\circ$ from the direction of travel). In addition to the images, we record all the boat sensor data: position from GPS, heading from compass, pitch and roll angle from IMU although they can be neglected and proximity from the laser range finder (not used beyond the on-board controller in this study).

This study includes data from 80 surveys from the second half of 2013 up to the end of 2015. This corresponds to potentially a bit more than 3'000'000 images collected over 80km of autonomous navigation.

The particularity of our dataset is that most of our images depicts natural scenes combining some water, trees and shrubberies at various distances, sometimes grass areas and/or far-away buildings, and sky. All of these elements are challenging for computer vision: the lake surface acts as a somewhat deformable mirror, sometimes very smooth and reflective and at other times not reflective at all due to wavelets. Additionally, flooding events means that the water line can move by up to 1m in some surveys. Trees are challenging for three reasons, first these ones do not always have leaves, second they are fractal self-similar structures, and last they are 3D semi-transparent structures whose appearance is very sensitive to view point, especially in winter. Finally, the sky varies with the weather and the sun position. Because we run the boat on the perimeter of the lake at different times of the day and as long as it is not raining (to avoid water drops on the camera dome), our images are also sometimes affected by sun-glare or very challenging dynamic

range requirements.

III. RELATED WORK

Mounting evidence suggests that the traditional approach to data association, i.e., using local image features, is unreliable in unstructured environments. It is more applicable the more structured an environment is. Point-based features can be associated well indoors, but special care has to be taken as they are applied in urban environments (e.g., street-view) [3], [4]. They lose representational power as the environment changes with night [5], rain [6], and shadows [7]. This means that in some natural environments, like lake shores, point-based feature matching is sporadic even among images from the same survey, and is unreliable between different surveys [8].

The lack of a dominant method for data association in outdoor environments has led to a number of new approaches. All of them function using some form of information beyond the capabilities of point-based features. Image sequence [9], [10], [11], [12], [13], image patch [14], [15], and whole image [16], [17] techniques are becoming increasingly dependable. There are, however, still shortcomings among them. A common limitation is robustness to changes in viewpoint among some approaches based on sequences or whole images. This may not be a factor in monitoring applications, however, since surveys are captured from similar trajectories; the viewpoint and the scale are relatively stable between images (see e.g., [18]).

In the recent years, deep neural networks have become very popular methods for solving both classification and regression problems because technical difficulties related to their training have been overcome. In the context of image analysis, CNNs have been around for several decades because they benefit from inherent regularities in images to constrain the trained architecture and their architecture regularize more general deep neural networks. In the recent years, deep CNNs have set multiple benchmarks for state-of-the-art performance on various machine learning tasks [19], [20]. Of particular interest for our study, deep CNNs have been successfully applied to place recognition [21] (a classification task), pose regression [22] and viewpoint estimation [23]. Finding the best neural network architecture for solving a given machine learning problem can be very challenging. In this study, we consider the CaffeNet architecture which is an implementation in Caffe[2] of the AlexNet CNN [19]. This network had state of the art classification accuracy on the ImageNet Large Scale Visual Recognition Challenge.

IV. METHODOLOGY

A. Data Pre-Processing

We consider two experimental setups: a classification task and a regression task. For the classification task, each image is labeled with the discretized pose of the robot. The pose consists in the position (from the GPS) and the heading (from the compass) of the robot. The position is discretized into a 350 m by 600 m grid, with 2.5 meters squares, centered over the lake. The heading is discretized in nonoverlapping

angular sectors of 10 degrees. This formulation has a total of $1'209'600$ unique labels.

A subset of the label space was observed in our dataset. Of the different labels that were observed, some were eliminated in order to obtain a balanced training set. That is, a class was only used if it was at least half as big as the largest class, which contained 1750 images. The training set consisted of 295 classes with approximately 1'000 images for a total of roughly 300'000 images. 5000 images were randomly selected for the test set.

The regression task used the same set of images and the labels were defined from the pose of the robot. An alternative was to use the position and the heading of the robot as labels, but this would imply to define a specific loss taking into account the angular nature of the heading. Although feasible in Caffe, this requires an in-depth modification of the library. Instead, we used a Euclidean loss and therefore defined the labels as a four-component vector with the position of the robot (from GPS) and the position of the point 10 meters away from the robot along the optical axis of the camera. One potential drawback of this approach is that the regression problem becomes more complicated than if we were to predict the position and heading since it requires the regressor to predict a specific location along the heading. However it turns out that despite this constraint, the regressor performed reasonably well (see section V-B). Finally, in order to ease the definition of the learning rate and to speed up convergence of learning, the labels to be regressed are normalized and centered.

B. CNN architecture and training

In this study, we used the AlexNet CNN[19] trained in Caffe[2]. The network consists in five convolution layers, five pooling layers, seven rectified linear unit layers, two normalization layers and three fully connected layers. Minor modifications were required to successfully train AlexNet. For the classification task, the size of the mini-batches is decreased to 32 samples and learning rate was decreased at a regular rate 10 times throughout training. The output layer of the fully connected part of the architecture uses a softmax transfer function to get a probability distribution over the labels and the loss is the cross-entropy classification loss.

For the regression task, a linear output transfer function is considered and the loss is the Euclidean loss. Experimentally, it was required to consider a lower learning rate than AlexNet which otherwise lead to a divergent loss. As we shall see in the result section, several strategies for setting the learning rate are considered. For both the regression and classification problems, the architecture is trained with the default CaffeNet settings, namely stochastic gradient descent, a momentum set to 0.9 and a weight decay to 0.0005.

C. Extracting season invariant representations

Being able to classify an image as belonging to one part of the lake with its viewpoint or to regress from it the pose of the boat are of interest by themselves. However, one of the objectives of the study was also to extract season

invariant representations in order to detect the changes of the lake shore. This part of the study was done using the network trained for the classification task. For every class of the selected dataset, a prototypical image was computed by averaging all the filter responses, at a given depth of the network, of all the images belonging to the considered class. A query image is then propagated through the network up to the depth where the prototypes have been computed, the responses of all the filters at that depth are then averaged and this representation is compared to the computed prototypes. The quality of a prototype image is assessed using the cosine similarity between the representation of the query image and the prototype. A new image is labeled by picking the class whose prototype has the largest similarity with the image's averaged representation.

V. RESULTS

A. Classification

Our best training on the aforementioned dataset was made with 300 000 iterations over mini-batches of size 32 (down from the default 256 for better accuracy), with a learning rate decreasing at a fixed rate ten times during the course of training. The full training took 5 days to finish on a Tesla K20C machine, and attained 70% accuracy on the test set, on a top-1 classification basis. Such results are satisfactory considering the similarities of natural scene images. It should be noted that classification was not the main goal, but a good classification will intuitively lead to better class representations.

We can thus extract the trained filters responses to see how images are processed by the network. The conv1 layer mostly learned edges, namely the skyline and the waterline. Conv2 detected foliage, and convolutional layers 3 through 5 contained low-level features. We tested prototype generation on all convolutional layers, as well as the last pooling layer pool5, shown in Fig. 3.

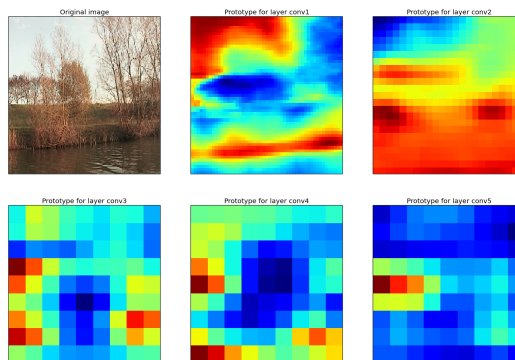


Fig. 3. Example of prototypes for a given image

We generated prototypes for all classes, and tested whether an image can be recognized only using its class prototypes. Testing over a thousand random images and measuring using a cosine distance shows that all layers can be used as suitable descriptors, with distances to the wrong prototypes being

indubitably larger than distances to the right prototype on average (see Table I, column 2 and 3 and Fig. 4)

Layer	Overall Ratio (full dataset)	Precision (%) top-20	Ratio (seen) (1 st half)	Ratio (unseen) (2 nd half)
conv1	1.76	43.9%	1.70	1.03
conv2	2.99	42.3%	0.93	0.96
conv3	1.86	34.8%	1.40	0.94
conv4	1.79	08.4%	1.16	1.00
conv5	1.68	57.3%	1.32	0.98
pool5	1.86	52.0%	1.42	0.95

Ratio of median distance of a random image to the wrong prototypes over median distance to the correct prototype. 2nd column refers to the ratio achieved using the complete dataset for training. 3rd column give the percentage of successful top-20 classification using the distance to the prototypes of the full dataset. 4th column is similar but using only half of the classes. 5th column evaluate the generalization performance by evaluating images from the classes not used for training.

TABLE I
CLASSIFICATION PERFORMANCE

We also trained the same network on half the classes, to test for generalization capabilities. We tested whether the network learned to reduce the dimensionality of an image into its season-invariant representation. In the case of the classes observed in the training set, the representation results are analogous to the full dataset. However, the internal features were not discriminative when applied on images from unseen classes (Table I, columns 4 and 5). The network learned to discriminate between its known classes, but it did not produce a season-invariant representation.

B. Regression

The objective of the regression task is to perform localization with performance comparable to an inexpensive GPS system. The training parameters were tuned to reach the best performance in the regression task.

The layers of the network were initialized using normal distributions. The initialization variances were set to 0.1 for Conv1 and 0.05 for Conv2, Conv3, Conv4 and Conv5. This made the weights big enough to propagate information through the network while still ensuring convergence. The initial learning rate and its evolution policy heavily depends on the loss function. For the regression task, very high values and a diverging behaviour were observed with the Euclidean loss at the beginning of the training. The learning rate and the weight decay on the last fully connected layer had to be set to 5e-04 and 2.5e-6, respectively, to avoid these effects. The learning rate policy was set to the “step” policy from Caffe and it was chosen to decrease the learning rate by half every 25’000 iterations. In our case this corresponds to the number of iterations required to observe the stabilization of the loss after each exponential decrease.

In this context, we tested three different approaches. The first one consists in training every convolution layer with the same learning rate. The second one consists in fine-tuning the learning rate of the Conv1 layer based on the results given by the first approach. The third one consists in loading the weights of the convolution layers from the classification

a) First approach b) Second approach c) Third approach

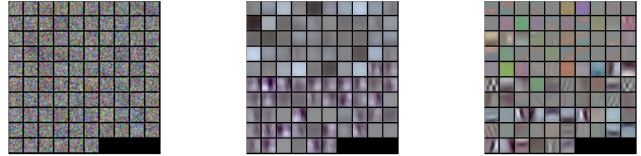


Fig. 6. Conv1 filters after 135 000 iterations.

training. In order to compare our results with those three approaches, we refer to the following loss function :

$$Loss = \frac{0.5}{scale^2} * \sum_{i=1, 4} (label_i - prediction_i)^2 \quad (1)$$

The first approach allowed to achieve an average error of 19.3 meters per label as it is shown in Fig. 5.a. The loss computed on the test dataset is plotted in green and the loss computed on the train dataset is plotted in blue. However after the training, the convolution layer filters did not exhibit any specific features (Fig. 6.a). Thus, it appears that the regression was only supported by the fully connected layers. Because our goal is to build a season-invariant representation of natural scenes, this approach did not reach our expectations.

The second approach led us to push further the difference of behavior between the fully-connected layers and the convolution layers. By doubling the learning rate and the weight decay for Conv1, we forced the convolution layers to contribute more to the regression. This method resulted in being more successful than the previous one. The best average error we achieved was 18.3 meters (Fig. 5.b). Some natural environment features can be identified in the convolution filters (Fig. 6.b).

The last approach used the same learning rate settings as the first approach. It was observed that the convolution weights decreased during the training and ended with a distribution similar to the second approach. The best average error achieved was 20.9 meters (Fig. 5.c). However the convolution filters retrieved exhibited different structures than the second approach (Fig. 6.c). Based on this result it can be concluded that the convolution weights learned during the classification training could not be reused for the regression task. The weights needed for this task required a smaller variance and the final model presented the worst results among the three approaches. Consequently, learning from normally distributed convolutional gains seems to be more efficient in the case of a regression.

The best results regarding the loss values on the train and test datasets were achieved with the second approach. After 150’000 iterations, the prediction errors are centered Gaussian-like distributions with a standard deviation of 11.8 meters on X and 21 meters on Y. The labels were displayed to be compared to the predicted positions (Fig. 7.a and b). The red and green dots represent the original labels and the blue and purple dots represent the predicted labels. The error

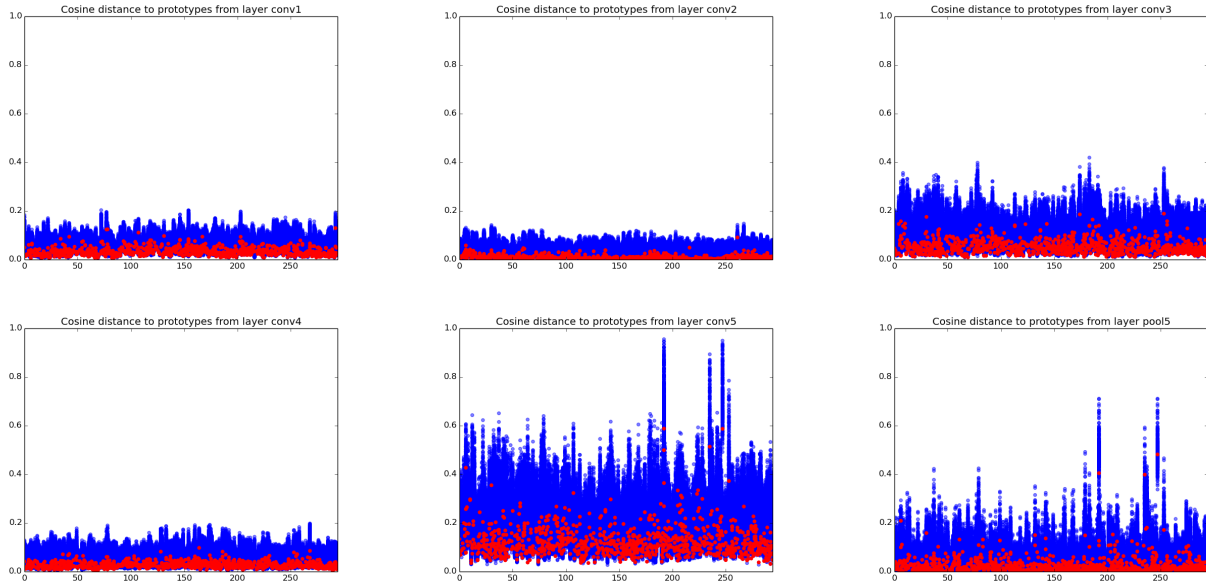


Fig. 4. Red dots represent the distance from a random class image to its class prototype, blue dots are distances to other class prototypes. Graphs show layers conv1 through pool5.

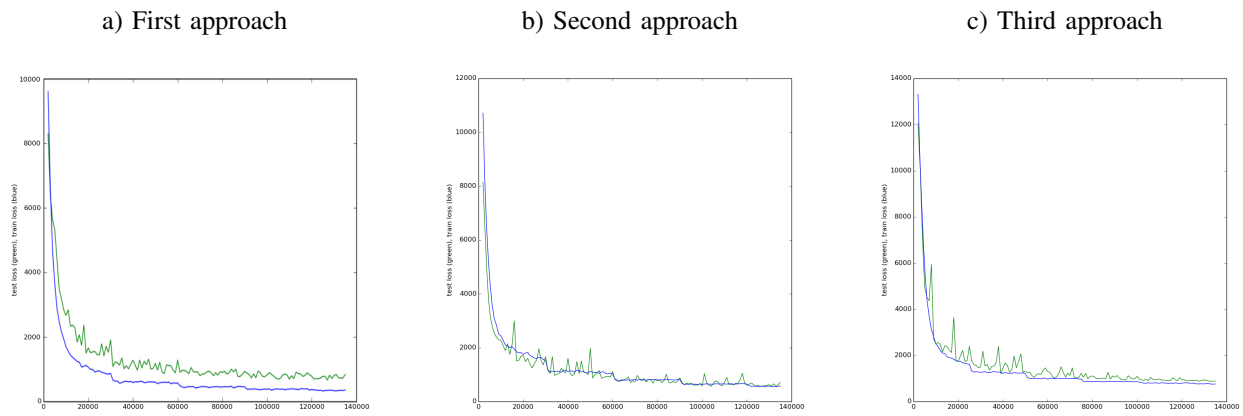


Fig. 5. Loss Function for 135 000 iterations.

between X and Y coordinates of the predicted and original labels was also plotted on Fig. 7.c and was found to be centered on the origin without significant bias.

VI. CONCLUSIONS

This paper evaluated the performance of off-the-shelf CNNs on a place recognition task and a pose prediction task for natural environments under large seasonal changes, in the context of a long-term autonomous monitoring problem. To this end, we presented an original dataset consisting in several million images taken at weekly interval on the shore of a small lake over two years.

The inconsistency of the appearance of the water and the sky, as well as the strong seasonal changes of vegetation and the weather-dependent lighting conditions proved to be manageable both for the classification task (70% precision) and for the pose regression task (20m standard deviation over 1km of shore line). However, it turned out that using

the standard network architecture did not result in learning generalizable features leading to a season-invariant representation of the environment. Because Caffe was too limited to explore this possibility within this study, a more general network architecture (as in [24]) would be appropriate.

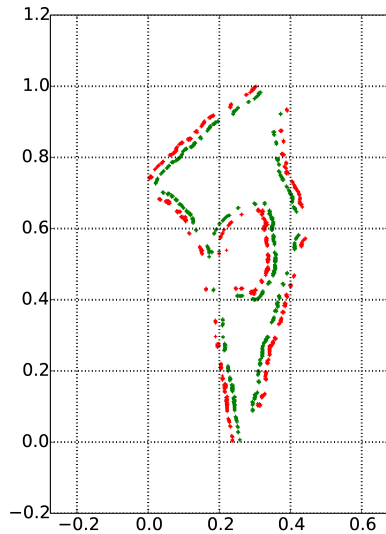
ACKNOWLEDGMENT

This work is supported by the Lorraine Region (France) and by the European Commission within the Flourish project, EC Grant agreement 644227.

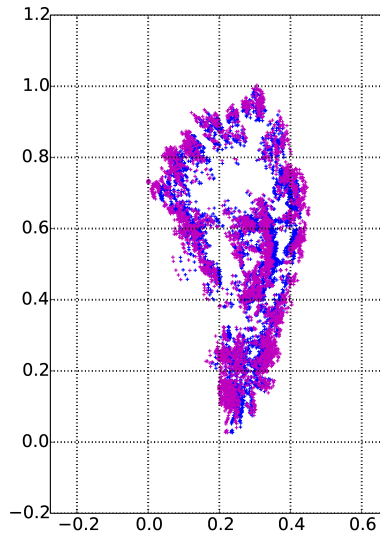
REFERENCES

- [1] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, August 2013.
- [2] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional Architecture for Fast Feature Embedding," *arXiv preprint arXiv:1408.5093*, 2014.
- [3] C. Beall and F. Dellaert, "Appearance-based localization across seasons in a Metric Map," in *6th PPNIV*, Chicago, USA, September 2014.

a) Map of original labels



b) Map of predicted labels



c) 2D error of the predicted label position

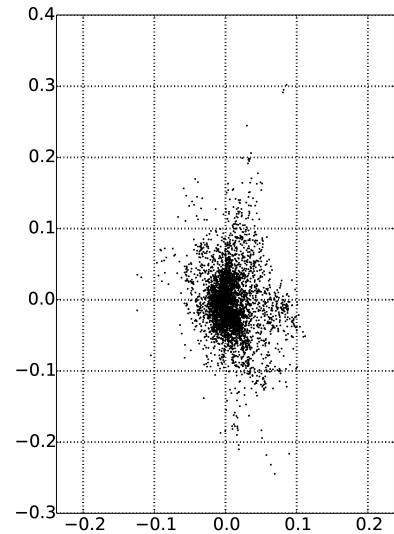


Fig. 7. Map of predicted positions and original labels.

- [4] E. Stumm, C. Mei, and S. Lacroix, "Probabilistic place recognition with covisibility maps," in *Intelligent Robots and Systems (IROS)*. IEEE, 2013, pp. 4158–4163.
- [5] P. Nelson, W. Churchill, I. Posner, and P. Newman, "From Dusk till Dawn: Localisation at Night using Artificial Light Sources," in *ICRA*, 2015.
- [6] A. Cord and N. Gimonet, "Detecting unfocused raindrops: In-vehicle multipurpose cameras," *Robotics & Automation Magazine, IEEE*, vol. 21, no. 1, pp. 49–56, 2014.
- [7] P. Corke, R. Paul, W. Churchill, and P. Newman, "Dealing with shadows: Capturing intrinsic scene appearance for image-based outdoor localization," in *IROS*. IEEE, 2013, pp. 2085–2092.
- [8] S. Griffith, P. Drews, and C. Pradaliere, "Towards autonomous lakeshore monitoring," in *International Symposium on Experimental Robotics (ISER)*, 2014.
- [9] M. J. Milford and G. F. Wyeth, "Seqslam: Visual route-based navigation for sunny summer days and stormy winter nights," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 1643–1649.
- [10] M. Cummins and P. Newman, "Fab-map: Probabilistic localization and mapping in the space of appearance," *The International Journal of Robotics Research*, vol. 27, no. 6, pp. 647–665, 2008.
- [11] M. J. Milford, G. F. Wyeth, and D. Rasser, "Ratslam: a hippocampal model for simultaneous localization and mapping," in *Robotics and Automation, 2004. Proceedings. ICRA'04. 2004 IEEE International Conference on*, vol. 1. IEEE, 2004, pp. 403–408.
- [12] W. Churchill and P. Newman, "Experience-based navigation for long-term localisation," *IJRR*, vol. 32, no. 14, pp. 1645–1661, 2013.
- [13] T. Naseer, M. Ruhnke, C. Stachniss, L. Spinello, and W. Burgard, "Robust visual slam across seasons," in *IROS*, 2015.
- [14] C. McManus, B. Upcroft, and P. Newman, "Scene signatures: Localized and point-less features for localization," in *RSS*, Berkeley, USA, July 2014.
- [15] N. Sunderhauf, S. Shirazi, A. Jacobson, F. Dayoub, E. Pepperell, B. Upcroft, and M. Milford, "Place recognition with ConvNet landmarks: Viewpoint-robust, condition-robust, training-free," *Proceedings of Robotics: Science and Systems XII*, 2015.
- [16] R. Arroyo, P. F. Alcantarilla, L. M. Bergasa, and E. Romera, "Towards life-long visual localization using an efficient matching of binary sequences from images," in *Robotics and Automation (ICRA), 2015 IEEE International Conference on*. IEEE, 2015, pp. 6328–6335.
- [17] P. Neubert, N. Sunderhauf, and P. Protzel, "Superpixel-based appearance change prediction for long-term navigation across seasons," *Robotics and Autonomous Systems*, vol. 69, pp. 15–27, 2015.
- [18] M. Milford, J. Firn, J. Beattie, A. Jacobson, E. Pepperell, E. Mason, M. Kimlin, and M. Dunbabin, "Automated sensory data alignment for environmental and epidermal change monitoring," in *Australasian Conference on Robotics and Automation 2014*. Australian Robotic and Automation Association, 2014, pp. 1–10.
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105.
- [20] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [21] N. Sunderhauf, F. Dayoub, S. Shirazi, B. Upcroft, and M. Milford, "On the Performance of ConvNet Features for Place Recognition," *CoRR*, vol. abs/1501.04158, 2015.
- [22] T. Pfister, K. Simonyan, J. Charles, and A. Zisserman, "Deep Convolutional Neural Networks for Efficient Pose Estimation in Gesture Videos," in *ACCV (1)*, ser. Lecture Notes in Computer Science, D. Cremers, I. D. ReidZ, H. Saito, and M.-H. Yang, Eds., vol. 9003. Springer, 2014, pp. 538–552.
- [23] H. Su, C. R. Qi, Y. Li, and L. J. Guibas, "Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views," in *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [24] A. Radford, L. Metz, and S. Chintala, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks," *arXiv preprint arXiv:1511.06434*, 2016.