

# Unsupervised Place Discovery for Visual Place Classification

Fei Xiaoxiao    Tanaka Kanji    Inamoto Kouya    Hao Guoqing

Univ. of FUKUI

3-9-1, bunkyo, fukui, fukui, Japan

e-mail [tnkknj@u-fukui.ac.jp](mailto:tnkknj@u-fukui.ac.jp)

## Abstract

*In this study, we explore the use of deep convolutional neural network (DCNN) in visual place classification for robotic mapping and localization. An open question is how to partition the robot's workspace into places so as to maximize the performance (e.g., accuracy, precision & recall) of potential DCNN classifiers. This is a chicken and egg problem: If we had a well-trained DCNN classifier, it is rather easy to partition the robot's workspace into places, but the training of a DCNN classifier requires a set of pre-defined place classes. In this study, we address this problem and present several strategies for unsupervised discovery of place classes ("time cue", "location cue", "time-appearance cue", and "location-appearance cue") and evaluate efficacy of the proposed methods using publicly available NCLT dataset.*

## 1 Introduction

Visual place classification (VPC) is a fundamental task in robotic mapping and localization [1]. In it, a mapper robot collects a set of training images with ground-truth viewpoint information, assigns a class label (place ID) to each image, and learns an environment map from the labeled training data. Then, a map user robot takes a visual image without viewpoint information, and classifies it into one of the learned place classes. In this paper, we are motivated by recent success of deep convolutional neural network (DCNN) [2] in various classification tasks [3], and explore the use of a DCNN classifier as an environment map.

An open question is how to partition the robot's workspace into places. This is an important problem as the definition of place classes strongly influences performance (e.g., accuracy, precision & recall) of a VPC task. Intuitively, each place class should be defined as a continuous region in the robot's workspace with similar DCNN features. The main difficulty is a chicken and egg problem: If we had a well-trained DCNN classifier, it is rather easy to partition the robot's workspace into place regions, but the training of a DCNN classifier requires a set of pre-defined place classes.

In this study, we address this problem and present several strategies. It is assumed that we are given a collection of visual images with ground-truth viewpoint as a guide, which can be independent from training and test data. The goal is to search for an effective partition of the workspace into places so as to maximize performance of potential DCNN classifiers. We pro-

pose to use three different types of information: time cue, location cue, and appearance cue that is available from the pre-trained DCNN. Then, we present four different strategies for workspace partitioning by combining them: "time cue", "location cue", "time-appearance cue", and "location-appearance cue". We then evaluate efficacy of the proposed methods using publicly available NCLT dataset in [4].

## 2 Approach

### 2.1 System Overview

Fig.1 illustrates an overview of our approach. We assume a typical supervised classification framework for VPC. The entire VPC framework consists of two phases: (1) training, and (2) testing. The training phase ("Train" in Fig.1) takes as input a set of labeled training images for each place class and trains a classifier that classifies an image into one of the pre-defined place classes. We represent a label by the place class ID. The testing phase ("Test" in Fig.1) takes as input a novel unseen image ("Query image" in Fig.1) and predicts its place class by using the trained classifier. Our use of DCNN for learning and testing follows a typical transfer learning [5], where the DCNN is pre-trained on Big data and then fine-tuned to adapt to the target domain [2]. Our experimental system is based on Alexnet pre-trained on the ImageNet LSVRC-2012 dataset and fine-tuning ("Fine-tuning" in Fig.1) on the relatively small training dataset ("Dataset" in Fig.1), which our algorithm creates from the NCLT dataset ("Build dataset" in Fig.1).

Although our approach is sufficiently general and applicable to various types of environments (e.g., indoor and outdoor) and sensor modalities, in experiments, we focus on North Campus Long-Term (NCLT) dataset in [4]. The NCLT dataset is a large scale, long-term autonomy dataset for robotics research collected on the University of Michigan's North Campus by a Segway robotic platform. The Segway is outfitted with a Ladybug3 omnidirectional camera, a Velodyne HDL-32E 3D lidar, two Hokuyo planar lidars, an inertial measurement unit (IMU), a single-axis fiber optic gyro (FOG), a consumer grade global positioning system (GPS), and a real-time kinematic (RTK) GPS. The data we used in our research is including image data and navigation data from NCLT dataset. The image data is from the front directed camera (camera#5) of the Ladybug3 omnidirectional camera. Fig.2 shows a bird's eye view of the experimental environment and an example robot trajectory.



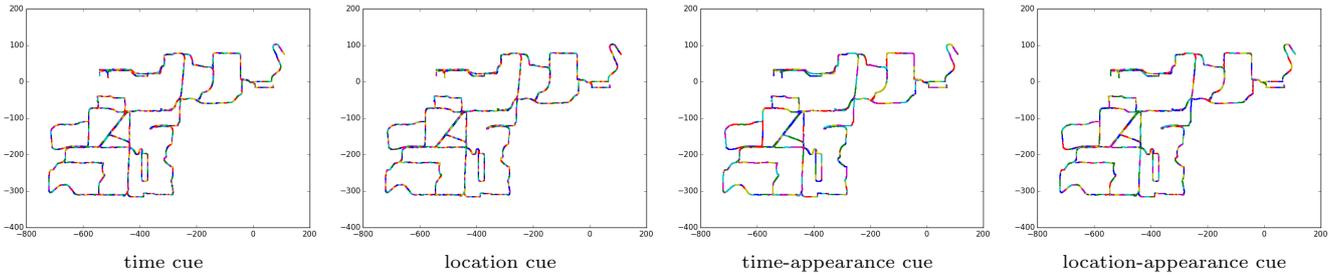


Figure 3. Example results for workspace partitioning.

of robotic mapping and localization, accuracy of VPC should reflect the cluster size  $|C(r_n)|$ , as the robot’s final goal is to localize a single test image rather than a cluster of test images. To address this issue, we also introduce a normalized version of the success rate (NSR) in the form:

$$p = \frac{1}{K} \sum_{n=1}^K \frac{|a_n \in r_n|}{|C(r_n)|}. \quad (3)$$

Note that  $|C(x_n)|$  serves as a regularizer to avoid trivial meaningless solutions in which all the training images belong to a single place class and other place classes are empty. Intuitively,  $p$  represents the probability of a given image being correctly classified into the ground truth place class.

## 2.4 Workspace Partitioning Strategies

We developed four different strategies for workspace partitioning. Fig.3 shows examples of place classes found by individual strategies. In the figure, different colors indicate different place classes. We can see that all the strategies create a set of clusters with similar cluster size, although performance difference between them is significant as shown in the experimental section, Section 3.

The first one is a simple time cue strategy. It partitions a sequence of images into classes by their time stamps or image IDs. Thus, the partitioning result is a set of clusters with  $K - 1$  intervals with approximately equal time length. This strategy is based on an observation that images with similar time stamps are expected to have visually similar appearance (i.e., DCNN features), as they are collected by a mapper robot that navigates through a continuous trajectory in the environment, and such clusters of images are expected to be a good training set for a DCNN classifier. Obviously, this simple strategy has many limitations. Particularly, it is not robust against variation of robot’s moving speed. In addition, it does not take advantage of any appearance features that are available from the pre-trained DCNN.

The second one is location cue strategy. This strategy is different from the time cue strategy only in that it partitions a sequence of images not by their time stamps but by their travel distance along the trajectory. Thus, the partitioning result is a set of clusters with  $K - 1$  intervals with approximately equal travel distances. Note that the location cue strategy does require the information of travel distance along the trajectory, which is readily available given the ground-truth viewpoint information. This strategy is robust

against variation of robot’s moving speed but still does not take advantage of appearance information from the pre-trained DCNN.

The third one is location-appearance cue strategy. The basic idea is to augment the location cue strategy by using the available pre-trained DCNN classifier as a guide. We use the 6-th layer from the pre-trained DCNN as the image representation, as it has shown excellent performance in image classification task in [6]. The workspace partitioning procedure is as follows. (1) First, images are represented by 4,096 dimensional 6-th layer features from the DCNN. (2) Second, they are fed to k-means clustering to obtain  $K$  image clusters. (3) Third, for each cluster, we perform the location cue strategy to partition the cluster into sub-clusters.

The fourth one is time-appearance cue strategy. The basic idea is to augment the time cue strategy by using the available pre-trained DCNN classifier as a guide. The basic concept of the augmentation is similar with that of the location-appearance cue strategy, but different in that we perform the time cue strategy (instead of the location cue strategy) at the step3 in the procedure.

## 3 Experiments

We evaluated the proposed framework for workspace partitioning on real VPC tasks. Four different strategies for workspace partitioning: time, location, location-appearance, and time-appearance, presented in the previous section were considered. For training data, we used dataset of “March 31st 2012”, in which the travel distance is 6.0km in total, the time is mid-day, the environmental condition is cloudy, no foliage and no snow. The training data consists of images and viewpoint information that is available from the NCLT dataset. For testing data, we used dataset of “Aug 4th 2012”, in which the travel distance is 5.5km, the time is morning, the environmental condition is sunny, foliage and no snow. The test data consists of images while the available viewpoint information is used for ground-truth prediction. Fig.4 shows examples of images belonging to individual clusters in the case of time-appearance strategy. Shown are random representative images for 600 clusters that are randomly sampled from the  $K$  clusters.

The training data provided by each strategy is fed to transfer learning (i.e., fine-tuning) of CNN that is pre-trained on the Bigdata (i.e., ImageNet LSVRC-2012 dataset). The classification function in the pre-trained CNN is a softmax classifier that computes the likelihood over 1,000 classes of the ImageNet dataset.

Table 1. Experimental result.

Strategy	#Class	top-5 SR	top-1 NSR	Loss
time	728	31.6%		3.97
location	728	33.6%		3.92
time-appearance	675	41.0%	0.938%	3.71
location-appearance	743	39.4%	1.664%	3.65



Figure 4. Visual image for each class.

To fine-tune the CNN, we change the number of the softmax classifiers at the top layer with the number of place classes. Then, the CNN parameters are fine-tuned on new training datasets. After the fine-tuning, we evaluate performance of CNN on the test set in terms of accuracy. In the experiment, we changed the softmax classifier with a new value which is equal to the classes of training datasets.

Table 1 shows the loss and accuracy of test datasets. We can see that the time-appearance cue strategy outperformed the other three strategies of time cue, location cue, and location-appearance cue. Besides, the location cue strategy outperformed the time cue strat-

egy. The reason may be that the time cue strategy does not consider the fact that the mapper robot moves with variable velocity and it often fails to partition the workspace into equal-size sub-regions. Overall, time-appearance and location-appearance strategies outperformed the other two. It could be said that appearance information from the pre-trained DCNN provides an effective cue to further improve the time cue strategy and the location cue strategy. Finally, the success rate of these strategy is sufficiently high considering the fact that the number of possible places is large, e.g., 675.

#### 4 Conclusions & Future Works

In this study, we explored the use of deep convolutional neural network (DCNN) in visual place classification (VPC) for robotic mapping and localization. It has been shown that the proposed strategies for workspace partitioning enabled effective discovery, learning and classification of place classes. Our research showed that we can use location features and appearance features to partition the robot’s workspace into places, to help better fine-tuning of the CNN, and to improve overall performance of visual place classification. Since our approach is simple and orthogonal to the choice of DCNN features and clustering algorithms, it is applicable to diverse VPC applications.

#### References

- [1] S. M. Lowry, N. Sünderhauf, P. Newman, J. J. Leonard, D. D. Cox, P. I. Corke, and M. J. Milford, “Visual place recognition: A survey,” *IEEE Trans. Robotics*, vol. 32, no. 1, pp. 1–19, 2016. [Online]. Available: <http://dx.doi.org/10.1109/TRO.2015.2496823>
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [3] D. Kumar, “Deep learning based place recognition for challenging environments,” 2016.
- [4] N. Carlevaris-Bianco, A. K. Ushani, and R. M. Eustice, “University of michigan north campus long-term vision and lidar dataset,” *The International Journal of Robotics Research*, 2015.
- [5] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, “How transferable are features in deep neural networks?” in *Advances in neural information processing systems*, 2014, pp. 3320–3328.
- [6] A. Babenko, A. Slesarev, A. Chigorin, and V. S. Lempitsky, “Neural codes for image retrieval,” in *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I*, 2014, pp. 584–599.