

Probabilistic Location Recognition using Reduced Feature Set

Fayin Li
fli@cs.gmu.edu

Jana Kořecká
kosecka@cs.gmu.edu

Technical Report GMU-CS-TR-2005-2

Abstract

The localization capability is central to basic navigation tasks and motivates development of various visual navigation systems. These systems can be used both as navigational aids for visually impaired or in the context of autonomous mobile systems. In this paper we describe a two stage approach for localization in indoor environments. In the first stage, the environment is partitioned into several locations, each characterized by a set of scale-invariant keypoints and their associated descriptors. In the second stage the keypoints of the query view are integrated probabilistically yielding an estimate of most likely location.

The emphasis of our approach in the environment model acquisition stage is on the selection of discriminative features, best suited for characterizing individual locations. The high recognition rate is maintained with only 10% of the originally detected features, yielding a substantial speedup in recognition. The ambiguities due to the self-similarity and dynamic changes in the environment are resolved by exploiting spatial relationships between locations captured by Hidden Markov Model. Once the most likely location is determined, the relative pose of the camera with respect to the reference view can be computed.

1 Introduction

The problem of localization is of interest in several applications including augmentation of human navigation capabilities and mobile robot localization. Two main variations of the localization problem have long been established in the robotics community and are known as global localization (also known as robot kidnapping problem) and pose maintenance.

In this paper we will focus on the global localization aspect and demonstrate how to solve it by means of location recognition. The approaches used to tackle the location recognition problem vary depending on means of acquiring the location database, representation of individual locations and methods of recognizing them. Although the location recognition problem shares many common aspects with general object recognition, it also differs in several important ways.

1.1 Related Work

Due to the different nature of the location recognition task, several representations of locations were proposed in the past. In the initial attempts to location recognition, the locations were represented by multi-dimensional color histograms [1]. Representations which enable coarser classification of indoor and outdoor scenes used responses to banks of filters with varying level of spatial integration include [18, 17]. In the context of mobile robot navigation robust versions of subspace methods, where individual views were represented as a points in the high-dimensional space have been applied in case omnidirectional cameras [2]. Given the subspace representation the pose of the camera is typically obtained by spline interpolation method, exploiting the continuity of the mapping between the object appearance and continuously changing viewpoint. Approaches which used local image descriptors for location representation have chosen affine or rotationally invariant features [10, 20] or local Fourier transforms of salient image regions [15]. Due to the locality of these image features, the recognition can naturally handle large amounts of clutter and occlusions. The sparser set of descriptors can be, in case of both global and local features, obtained by principal component analysis or various clustering techniques.

Several instances of pose maintenance and acquisition of metric environment models have been successfully solved in smaller scale environments [3, 14]. The applicability of these methods to large dynamically changing environment poses additional challenges and calls for alternative models.

1.2 Approach Overview

We propose to tackle the location recognition and localization problem by using a model of the environment represented by a set of locations and spatial relationships between them. Each location is represented by a set of views and their associated local scale invariant features. We present a novel technique for identifying most discriminative features for individual locations reducing the feature database to 10% of its original size, without forgoing the recognition accuracy. An associated likelihood model characterizing each location is then used in the Hidden Markov Model framework which enables us to resolve misclassification due to the self-similarity and dynamic changes in the environment. Once the most likely location is determined we can compute the relative pose of the camera, with respect to the reference view. We will report on the localization experiments in indoor environment with 18 locations and discuss current implementation efforts toward real-time demonstration of the proposed system.

2 Location Representation

As a starting point of our method, we use the environment model obtained in the exploration stage. Given a temporally sub-sampled sequence acquired during the exploration, the sequence is partitioned into $N = 18$ different locations. The locations in our model correspond to hallways, sections of corridors and meeting rooms approached at different headings. The initial model was obtained by a mobile robot, which was guided through the environment. The path of the exploration route and labels associated with the individual locations are in Figure 1. The number of views per location vary between 5 to 20 depending on the appearance variation within the location. The transitions between the locations occur either at places where navigation decisions have to be made or when the appearance of the environment changes suddenly. The images were taken approximately every 2-3 meters. The representative views of some locations in Figure 2 demonstrate the variability of our dataset. More details about the model acquisition stage can be found in [8].

Individual locations are represented by scale-invariant (SIFT) keypoints described in [10]. The SIFT features

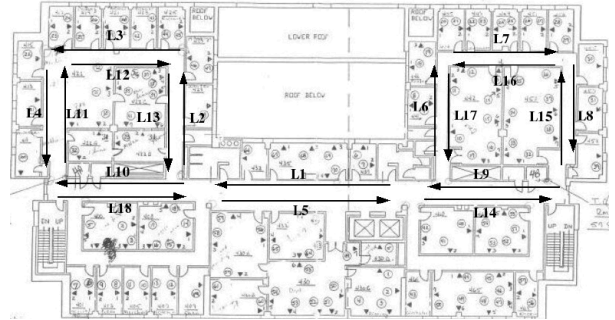


Figure 1: The map of the fourth floor of our building. The arrows correspond to the heading of the robot and the labels represent individual locations.



Figure 2: Examples of representative views of 12 out of 18 locations.

represent distinguishable image locations, which are stable across variations in scale. Each feature is endowed with a 128 dimensional descriptor, which captures the orientation information of local image region centered at each keypoint, is rotationally invariant and has been shown to be robust with respect to large variations in viewpoint and scale.

Figure 3a and 3b show the features detected in one of the representative views of locations 1 and 3. The number of features detected in each image varies from hundreds to thousands as shown in Figure 3c for the training data set. Using the matching scheme proposed in [10], the reliability of the match is measured by the ratio between the Euclidean distance to the closest neighbor and that to the second-closest neighbor. Figure 3d shows the number of matched features between consecutive views of the training sequence. Despite the large overlap between consecutive views, only a small number of features detected in consecutive views are matched. Some features have better capability to handle variations in scale and viewpoint and match stably in several different views of the same location. Selecting such features can keep the environment model more compact and save the computational cost for future localization. This ob-

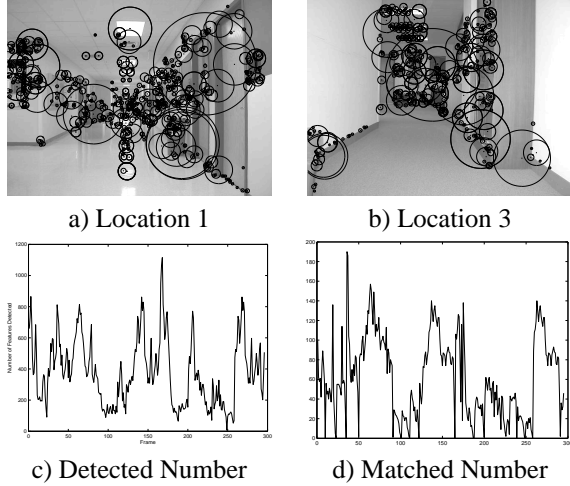


Figure 3: a) and b) 700 and 400 SIFT features detected in representative views of location 1 and location 3. c) Number of detected in the training sequence. d) Number of features matched between consecutive views in the training sequence.

servation brings to forefront the issue of feature selection in the model acquisition stage.

2.1 Feature Selection

Previously proposed techniques for reducing the feature pool include k-means clustering, greedy techniques or boosting [19, 4, 16, 7, 9]. Questions focusing on the model compactness [5] as well as trade-off between the complexity of features and the complexity of classifiers were explored in [11]. In [11] authors selected features with a greedy process, where only the features, which can increase the information content of the feature set with respect to the object were selected. In [5] authors estimate the posterior of each feature with respect to each object and the Shannon entropy is used to select the discriminative regions. Our method for feature selection is similar to the method proposed by [5] but with a different selection criterion. Suppose location L_i has N_i training images with total K_i of detected features $G_i = \{g_k^i\}_{k=1 \dots K_i}$. To obtain the information content of each feature g_k^i with respect to location identification, we need to estimate the posterior probability $P(L_l|g_k^i)$, $l = 1 \dots N$. The posterior probability at feature g_k^i is estimated using only features g_j inside a Parzen window [9] of a local neighborhood $Z = \{g_j | \|g_k^i - g_j\| \leq \epsilon, j = 1 \dots z\}$, where ϵ determines the size of the window. We weight the contribution of specific feature g_j^l in Z -labeled by location L_l - that should increase the posterior estimate $P(L_l|g_k^i)$

by a Gaussian kernel function $N(\mu, \sigma)$ in order to favor the features with smaller distance to the feature g_k^i , with $\mu = g_k^i$ and $\sigma = \epsilon/2.5$. Then the posterior probability at feature g_k^i is estimated as

$$P(L_l|g_k^i) \propto \sum_{g_j^l \in Z} \exp\left(-\frac{\|g_j^l - g_k^i\|^2}{2\sigma^2}\right). \quad (1)$$

There are two different ways how to determine the local neighborhood for each feature. One, is setting ϵ to some predefined threshold T , which may be estimated from the pairwise distance distribution of all features in the training images. In this case, the number z of elements in the set Z varies for each feature g_k^i , while the σ is constant for all features. Some features, however have a large number of features in the neighborhood while some have very few.

Since the threshold T is not easily determined, we choose the second approach where the number of elements in the feature neighbourhood is kept fixed. In this case the Parzen window size σ varies for each feature g_k^i , where σ is proportional to the largest distance ϵ from g_k^i to Z . We choose $z = 200$ and make the Parzen window adaptive.

After the posterior probabilities of all features are obtained, we could proceed by calculating the information entropy of each feature and use it in the selection process. In our data set, however, each location has different number of training images and the number of features detected in each image varies largely as it can be seen in Figure 3c. If one location l has few features detected and the Parzen window size is not accurate, even though a good feature g_k^l has a large posterior probability $P(L_l|g_k^l)$, it may also have a large posterior probability with respect to another class $P(L_m|g_k^l)$, because location m has large number of detected features. The posterior probability estimates will hence be biased and the entropy will not successfully capture the right information content.

Due to this reason we use directly the estimate of $P(L_l|g_k^l)$ in each image for ranking the features based on their discrimination capability. The number of features we keep is specified as a percentage η of detected features. For a training image from location L_i with M_i features detected, only top $\max(M_i\eta, N_0)$ features are selected based on their posterior $P(L_l|g_k^i)$ ranking. N_0 is the minimal number of features selected to avoid discarding too many features from the images of locations with few detected features. Because we select the features based on the rank of their posterior, the feature with high rank of posterior with respect to location L_i may also have large posterior with respect to another location L_j . That feature can be shared by several locations and can distinguish them from others. At the same time it can

introduce ambiguities in discriminating locations. Currently, instead of fully exploiting the shared features, as it was done in the context of hierarchical object detection approach proposed by [16], we properly account for their contribution in the feature matching stage.

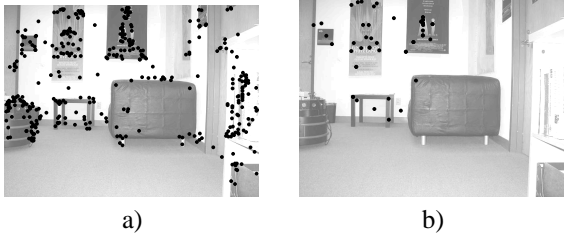


Figure 4: a) The total of 480 features detected and b) 50 informative features selected by our method.

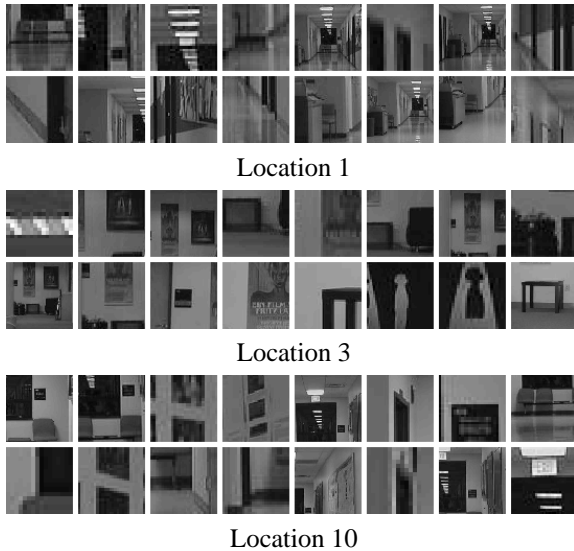


Figure 5: The top 16 features selected for three of the locations.

In our data set, there are 296 training images for 18 locations with 112,705 detected features. We chose $\eta = 10\%$ and $N_0 = 50$ to select features from each image. Each feature has the location (x, y) , a scale s , an orientation and a 128 dimensional descriptor. Figure 4 shows the total number of detected SIFT features and the informative features selected in the image by our method. The features belonging to posters have good discrimination capability. Figure 5 shows the top 16 selected features for different locations. Each feature is cropped from the training image centered at (x, y) with radius $r = 6 \times s$, where s is the scale of the SIFT feature. The patches are normalized to 64×64 .

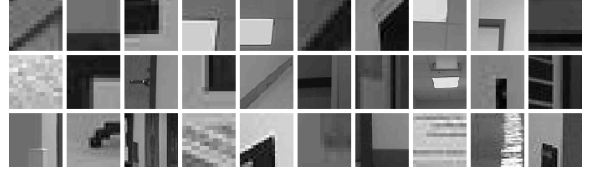


Figure 6: The features discarded in the selection.

Note that the selected features often have large scale and capture the global information about individual locations. For example for Location 1, which has large depth variation some the selected discriminative features are centered in the middle of the corridor and entail the entire view of the corridor. This also demonstrates that the suitable choice of the representation for a location varies largely and for certain locations the global descriptors are indeed highly discriminative. Figure 6 shows some examples of discarded features. For example features belonging to the ceiling lights are discarded except when the feature contains multiple lights.

3 Reduced Feature Set Matching

In order to demonstrate the feature selection process is effective, we compare the performance of location recognition using the reduced feature set, with the standard voting approach which uses all features. In this experiment the i -th location is represented by a number of representative views $\{I_n^i\}$ and their associated original and reduced SIFT feature sets $\{g_k(I_j^i)\}$ and $\{\tilde{g}_k(I_j^i)\}$. For a new query image Q and its associated features $\{g_k^Q\}$, a set of matches between Q and each model view I_j^i is determined by matching each feature in $\{g_k^Q\}$ against the model database features and choosing the nearest neighbor based on the Euclidean distance between two descriptors. Only the point matches whose nearest neighbor is at least 0.6 times closer than the second nearest neighbor are considered. More detailed justification behind the matching strategy and the choice of the threshold can be found in [10]. The model view I_j^i with the highest number of matched keypoints with Q is considered to be the correct result. To evaluate the proposed feature selection mechanism we compare the recognition performance using the model database of all detected features $\{g_k(I_j^i)\}$ and the reduced feature set $\{\tilde{g}_k(I_j^i)\}$.

Table 1 shows the recognition rates for the training sequence and two test sequences using the original and reduced feature set, respectively. The results are reported on the training sequence of 296 images from which the model views and features were selected and two test se-

sequence (# of frames)	original set	reduced set
No.1 (296)	100.0%	98.0%
No.2 (134)	82.1%	79.9%
No.3 (130)	83.1%	73.9%

Table 1: Recognition performance in term of percentage of correct localization using voting scheme.

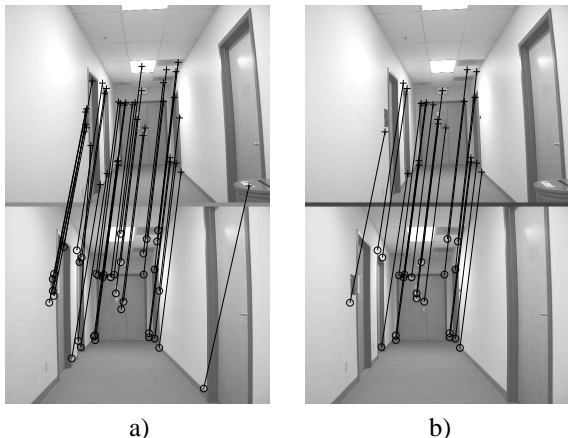


Figure 7: a) Matching using all detected features and b) using the selected features only. Top row: query image, bottom row: representative view.

quences of 134 and 130 images. The two test sequences, taken at different days and times of day, exhibit larger deviation from the path traversed during the training and several locations underwent dynamic changes which changed their appearance. In most cases, the matched features from the original set, were well preserved in the reduced feature set as shown in Figure 7. The selected features have enough discriminant ability to distinguish the locations. In few instances the reduced features yields correct recognition, while the wrong decision was made using the original feature set (see Figure 8). Relatively poor performance on the test sequences was due to several changes in the environment between the training and testing stage as demonstrated in Figure 11. Most SIFT features belong to objects some of which are not inherent to particular locations. In the next section we describe how to replace the voting scheme by a simple probabilistic model and propose how to resolve the remaining issues by explicitly modeling spatial neighborhood relationships between individual locations.

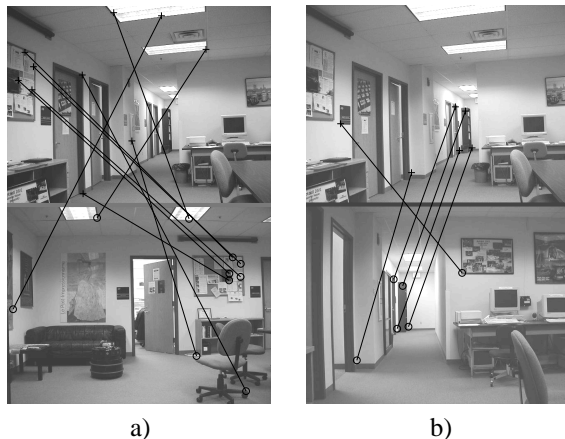


Figure 8: Selected features yield better matches than original feature set. a) the misclassification using original set b) correct recognition with one false match using the reduced set.

4 Probabilistic Location Recognition

The results of voting approach described in the previous section demonstrate that the feature selection is very effective. The misclassified locations are often due to the self-similarity of the environment (e.g. similarity of the appearance of corridors or hallways belonging to different locations), large changes in the pose between the query view and model views or dynamic environment changes. The classification performance can be improved by either exploiting more elaborate recognition scheme or additional information about the environment which would help to reduce the ambiguities due to the self-similarity of the environment. In the following section we demonstrate how to improve the classification by formulating the location recognition probabilistically and by exploiting the spatial relationships between the locations modeled by Hidden Markov Model.

In the voting framework, once the distance between two feature descriptors is within a specified threshold they are considered matched. The location with largest number of matched features is then declared to be the correct classification. The probabilistic formulation of the classification prediction entails computation of the posterior probability $P(L_l | \{g_k^Q\})$ of each location given the selected features from the query view. Such computation requires likelihood model for the matched features, which explicitly accounts for the quality of individual matches and hence is expected to be superior to the simple voting approach. The likelihood model can then be naturally used in the Hidden Markov Model (HMM)

to achieve more reliable and robust system.

When local descriptors are used as observations, several models of class posteriors have been proposed in the context of probabilistic approaches to object recognition [12, 13]. The proposed likelihood models account for the feature density and spatial relationships between features and have been shown to improve overall recognition rate. In the context of global image descriptors the locations were modeled in terms of Gaussian mixture models proposed in [17]. Those approaches have very complex parametric model and need large number of training examples to learn parameters. Furthermore the location recognition problem is notably simpler than the object recognition problem due to the background clutter not being so prominent¹. We propose a non-parametric method to estimate the $P(L_l|\{g_k^Q\})$ from training data directly without modeling the decision function. The essential features of this probabilistic method, which we describe next, is the selection of relevant features in the matching stage and integration of the evidence they provide for individual locations through a strangeness measure. It is a probabilistic version of voting approach.

As the result of the feature selection stage the features from the representative views of location i are joined to form a model of that location denoted by $\tilde{G}_i = \{\tilde{g}_k^i\}$. Given the query image Q with detected features $\{g_k^Q\}$ in order to determine most probable location, we need to compute the posterior probabilities $P(L_l|Q) = P(L_l|\{g_k^Q\})$ for $l = 1 \dots N$. Similarly as in the model building stage, many features in $\{g_k^Q\}$ are not informative and have no evidence for classification label. They may confuse the prediction if such features are considered during prediction, especially when cluttered background are present. We need to select *good* and relevant features from $\{g_k^Q\}$ for estimation of the posterior. The selection criterion not only gives the number of matched feature, but also yields the confidence of the match. The procedure is based on the hypothesis test.

Given a set of features in the query image $\{g_k^Q\}$, we first define the so called strangeness parameter α_k^i , which characterizes the discrimination capability of k -th feature, with respect to i -th location

$$\alpha_k^i = \frac{\min_{g_j \in \tilde{G}_i} (\|g_k^Q - g_j\|)}{\min_{g_j \notin \tilde{G}_i} (\|g_k^Q - g_j\|)}. \quad (2)$$

α_k^i is the ratio of minimal intra-distance within the class and minimal inter-distance to features from other classes. If α_k^i is greater than 1, the feature g_k^Q is not contributing to classification of Q as label L_i . The

¹The probabilistic models used in the object recognition, must also account for the fact that large number of detected features comes from background and not the object.

Sequence (# of frames)	Maximal Likelihood	HMM
No.1(296)	99.0%	100.0%
No.2(134)	85.8%	95.5%
No.3(130)	80.8%	95.4%

Table 2: Recognition Performance in term of percentage of correct localization based on α -values.

smaller the α_k^i is, the more discriminative is the feature for the purpose of classifying Q as i -th location. If $g_j^* = \operatorname{argmin}_{g_j \in \tilde{G}_i} (\|g_k^Q - g_j\|)$ is a shared feature among a set of locations S , the strangeness is very close to 1. Since we want the shared features to be considered, in the selection stage the strangeness of the shared feature g_k^Q with putative label l is re-computed as follows

$$\alpha_k^l = \frac{\min_{g_j \in \tilde{G}_l} (\|g_k^Q - g_j\|)}{\min_{g_j \notin \tilde{G}_l \wedge g_j \notin \tilde{G}_{i \in S}} (\|g_k^Q - g_j\|)}. \quad (3)$$

Hence in the case of shared features we do not consider the inter-distance from the features in shared locations $t \in S$ in this special case. The shared features can help to distinguish the location subset S from other locations but are useless for discriminating the locations in subset S . The computation of α -values has the same computational complexity as the nearest neighbor ratio computation in the standard voting scheme.

For the computation of the likelihood $P(\{g_k^Q\}|L_l)$ we select only top R features from the query image, ranked by their strangeness, under current hypothesis test. Note that don't consider the features who have strangeness measure $\alpha_k^l \geq 1$. The likelihood of feature g_k^Q has the putative label L_l is defined as

$$P(g_k^Q|L_l) = P(\alpha_k^l|L_l) = \exp(-\frac{\alpha_k^l{}^2}{2\sigma^2}). \quad (4)$$

Since we do not know how many features belong to location L_l among top R features, we need to integrate the evidence over all possible hypotheses. A hypothesis in our case indicates that a subset $h_j \neq \emptyset$ of top R features is classified as location L_l . Assuming the selected features are independent, we can now compute the probability of a single hypothesis h_j conditioned on location L_l

$$P(h_j|L_l) = \prod_m P(\alpha_m^l|L_l) \prod_n (1 - P(\alpha_n^l|L_l)). \quad (5)$$

Index m ranges over features which belong to location l with certain probability, n ranges over features which do not belong to location L_l , where $m + n = R$ is the number of selected features, i.e. the length of a hypothesis.

Then the probability $P(L_l|\{g_k^Q\})$ can be computed as

$$P(\{g_k^Q\}|L_l) = \sum_{h_j} \prod_m P(\alpha_m^l|L_l) \prod_n (1 - P(\alpha_n^l|L_l)). \quad (6)$$

It can be simplified to

$$P(\{g_k^Q\}|L_l) = 1 - \prod_{k=1}^R (1 - \exp(-\frac{\alpha_k^l{}^2}{2\sigma^2})). \quad (7)$$

In our experiments we use $\sigma = 1/3$. Assuming the location prior is uniform and $R = 10$, we tested the training and two test sequences again using the above maximum likelihood criterion. The recognition rates are shown in Table 2. The performance is very close to the one using the original training feature set while reducing the matching computational cost by about 90%. If we consider the informative factor of each training feature and can reliably estimate their likelihood, we can estimate the posterior by

$$P(L_l|\{g_k^Q\}) \propto P(L_l) (1 - \prod_{k=1}^R (1 - P(g_k^*|L_l) \times \exp(-\frac{\alpha_k^l{}^2}{2\sigma^2}))), \quad (8)$$

where $g_k^* = \operatorname{argmin}_{g_j \in G_l} (\|g_k^Q - g_j\|)$ is the feature from location l matched with g_k^Q .

In our method, the minimal number of features selected N_0 in each representative view is an important parameter to make the feature selection reliable. We vary N_0 from 20 to 70 and test the recognition performance again. The results are shown in Figure 9. It demonstrates that the feature selection process is very effective. The recognition rates will be almost constant after N_0 is greater than 60, which means using more features helps little in increasing the performance.

4.1 Exploiting Neighborhood Relationships

We propose further to deal with the dynamic changes in the environment by incorporating additional knowledge about neighborhood relationships between individual locations. The rationale behind this choice is, that despite the presence of ambiguities in recognition of individual views the temporal context should be instrumental in resolving them. The use of temporal context is motivated by the work of [17] which addresses the place recognition problem in the context of wearable computing application. The temporal context is determined by spatial relationships between individual locations and is modeled by a Hidden Markov Model (HMM). In this

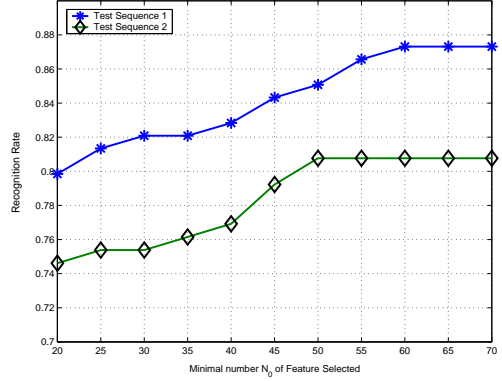


Figure 9: The recognition rates with different minimal number N_0 of feature selected.

model the states correspond to individual locations and the transition function determines the probability of transition from one state to another. We have already explored the use of this representation using original feature set and slightly different likelihood function [8]. Since the locations cannot be observed directly, each location is characterized by the location observation likelihood $P(o_t|L_t = L_l), l = 1 \dots N$ at time t during the exploration. The most likely location is at each instance of time obtained by maximizing the conditional probability $P(L_t = L_l|o_{1:t})$ of being at time t and location L_l given the available observations up to time t . The location likelihood can be estimated recursively using the following formula

$$P(L_t = L_l|o_{1:t}) \propto P(o_t|L_t = L_l)P(L_t = L_l|o_{1:t-1}) \quad (9)$$

where $P(o_t|L_t = L_l)$ is the observation likelihood, characterizing how likely is the observation o_t at time t to come from location L_l . The conditional probability $P(o_t|L_t = L_l)$ that a query image Q_t at time t characterized by an observation $\{g_k^{Q_t}\}$ comes from the location l is simply the likelihood $P(\{g_k^{Q_t}\}|L_t = L_l)$ introduced in Equation 4

$$P(o_t|L_t = L_l) \propto 1 - \prod_{k=1}^R (1 - \exp(-\frac{\alpha_k^l{}^2}{2\sigma^2})). \quad (10)$$

The second term of equation (9) can be further decomposed to explicitly incorporate the location neighborhood relationships

$$P(L_t = L_i|o_{1:t-1}) = \sum_{j=1}^N A(i, j)P(L_{t-1} = L_j|o_{1:t-1}), \quad (11)$$

where N is the total number of locations and A is a $N \times N$ matrix, where $A(i, j) = P(L_t = L_i|L_{t-1} = L_j)$

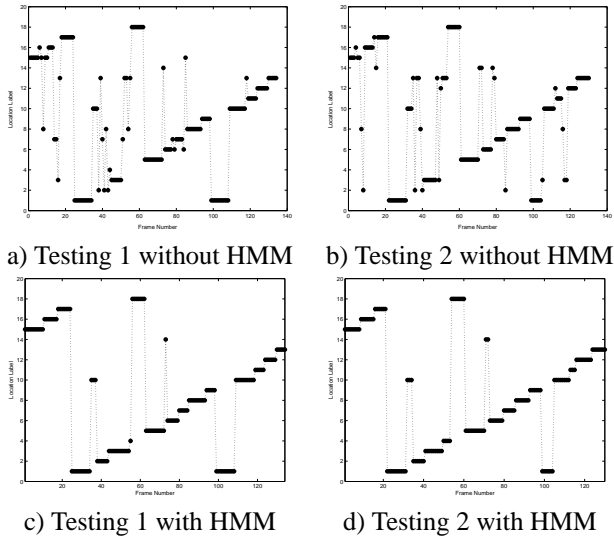


Figure 10: Classification results with for Test Sequence 1 and Sequence 2 with (bottom row) and without (top row) considering the spatial relationships modeled by HMM. The black circles correspond to the location labels assigned to individual frames of the video sequence.

is the probability of two locations being adjacent. In the presence of a transition between two locations the corresponding entry of A was assigned a unit value and in the final stage all the rows of the matrix were normalized.

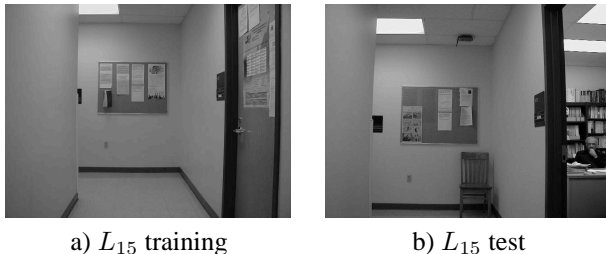


Figure 11: Appearance changes of location 15 between training and testing. There is no chair and the door is closed in the training view. The posters on the board are different between training and testing.

The results of location recognition employing this model are in Figure 10. For each frame of two test sequences, Figure 10 shows the location label which had the highest probability. The recognition rate with HMM for sequence 1 was 95.5% and for sequence 2 it was 95.4%. While in both cases some images were misclassified the overall recognition rates are a great improvement compared to the rates of single view location recognition. The dynamic changes, which make the

single view recognition fail are resolved successfully using HMM model. Figure 11 shows the example of dynamic changes. Despite some classification errors in test sequences, the order of visited locations was correctly determined. For test sequence 2, where we exhibited some intentional deviations between the path taken during training and testing, the classification of frames 69-70 as location 14 is incorrect (Figure 10d). The effect of HMM model can be examined by making all the probabilities in the transition matrix A uniform and essentially neglecting the knowledge of location neighborhood relationships. For comparison this is depicted in Figure 10a and 10b. Once the most likely location has been determined, we can estimate the relative pose of the camera with respect to the most likely representative view. This can be done by exploiting geometric relationship between two views captured by epipolar geometry. The detailed description of this stage in the context of the proposed application can be found in our earlier work [8].

4.2 Implementation and Experiments

The proposed approach for location recognition and localization has been tested in indoors environments and can operate both with or without the knowledge of the spatial relationships between the locations. Note that the recognition rate reported in Table 2 is relatively high even in the absence of HMM. The prerequisite of the approach is an off-line acquisition of the location database and the feature selection stage. In the testing stage the images can be acquired with a camera phone or camera equipped PDA, with the matching and recognition done at the location database server. The feature detection and matching, which is in our case optimized by making the model compact, has been already demonstrated in real-time in [6]. The SIFT feature extraction takes on average 150 ms for 640 x 480 image. Since the number of features in our model is significantly reduced, we expect the matching and recognition time to be superior to the previously reported results [6]. Currently the algorithm and model we proposed can recognize more than 4 frames per second without any special purpose optimization. Given the current implementation, an alternative mode of operation is to use the mobile robot platform as a guide to visually impaired person, where all the computing can be done on-board of the mobile platform. The advantage of this type of system is the fact that it can be integrated with additional guidance capabilities, such as obstacle avoidance and path finding.

We would also like to point out that the representation of locations which we proposed corresponds to semantically meaningful entities (e.g. corridors, hallways, conference rooms) or indoors environments. This may pose additional advantage in terms of facilitating better inter-

face between the PDA and visually impaired person.

5 Summary and Conclusions

We have demonstrated an approach for location recognition in indoor office like environments. The model of the environment is partitioned to individual locations and neighborhood relationships between them in the exploration stage. The individual locations were represented by SIFT features and location recognition was approached by feature matching between query and model views. We have presented a novel feature selection strategy, which exploited local information content and discriminability of the individual features and their associated descriptors. We have shown that by reducing the feature pool to 10% of the original size, we can achieve comparable performance to methods which use the original feature set. Further improvements were demonstrated by interpreting the quality of the features matches probabilistically and by endowing the environment with the HMM structure which exploits spatial relationships between locations. We are currently evaluating the effectiveness of the feature selection strategy in the context of other object and category recognition data sets.

Acknowledgments

The authors would like to thank D. Lowe for making available the code for detection of SIFT features. This work is supported by NSF grants IIS-0118732 and IIS-0347774.

References

- [1] H. Aoki, B. Schiele, and A. Pentland. Recognizing places using image sequences. In *Conference on Perceptual User Interfaces*, San Francisco, November 1998.
- [2] M. Artac, M. Jogan, and A. Leonardis. Mobile robot localization using an incremental eigenspace model. In *IEEE Conference of Robotics and Automation*, pages 1025 – 1030, 2002.
- [3] A. Davidson and D. Murray. Simultaneous localization and map building using active vision. *IEEE Transactions on PAMI*, 24(7):865–880, 2002.
- [4] R. Fergus, P. Perona, and A. Zisserman. Object class recognition by unsupervised learning. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 264–271, 2003.
- [5] G. Fritz, L. Paletta, and H. Bischof. Object recognition using local information content. In *ICPR*, 2004.
- [6] I. Gordon and D. Lowe. Scene modelling, recognition and tracking with invariant image features. In *International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 110–119, 2004.
- [7] S. Helmer and D. Lowe. Object class recognition with many local features. In *Workshop on Generative Model based Vision, CVPR*, 2004.
- [8] J. Kosecka and X. Yang. Location recognition and global localization based on scale-invariant keypoints. In *Workshop on statistical learning in vision, ECCV*, 2004.
- [9] N. Kwak and C. Choi. Input feature selection by mutual information based on parzen window. *IEEE Transactions on PAMI*, 24(12):1667–1671, 2002.
- [10] D. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 2004.
- [11] M. V. Naquet and S. Ullman. Object recognition with informative features and linear. In *ICCV, Nice, France*, 2003.
- [12] A. Pope and D. Lowe. Probabilistic models of appearance for 3-d object recognition. *International Journal of Computer Vision*, 40(2):149–167, 2000.
- [13] C. Schmid. A structured probabilistic model for recognition. In *Proceedings of CVPR, Kauai, Hawaii*, pages 485–490, 1999.
- [14] S. Se, D. Lowe, and J. Little. Global localization using distinctive visual features. In *Proc. of International Conference on Robots and Systems*, pages 153–158, 2002.
- [15] R. Sims and G. Dudek. Learning environmental features for pose estimation. *Image and Vision Computing*, 19(11):733–739, 2001.
- [16] A. Torralba, K. Murphy, and W. Freeman. Sharing features: efficient boosting procedures for multiclass object detection. In *International Conference on Computer Vision and Pattern Recognition*, 2004.
- [17] A. Torralba, K. Murphy, W. Freeman, and M. Rubin. Context-based vision system for place and object recognition. In *International Conference on Computer Vision*, 2003.
- [18] A. Torralba and P. Sinha. Recognizing indoor scenes. *MIT AI Memo*, 2001.
- [19] P. Viola and M. Jones. Robust real-time object detection. *International Journal of Computer Vision*, 57(2):137–154, 2004.
- [20] J. Wolf, W. Burgard, and H. Burkhardt. Using and image retrieval system for vision-based mobile robot localization. In *Proc. of the International Conference on Image and Video Retrieval (CIVR)*, 2003.