

A Visual Exploration Algorithm using Semantic Cues that Constructs Image based Hybrid Maps

Aravindhan K Krishnan and K Madhava Krishna

Abstract—A vision based exploration algorithm that invokes semantic cues for constructing a hybrid map of images - a combination of semantic and topological maps is presented in this paper. At the top level the map is a graph of semantic constructs. Each node in the graph is a semantic construct or label such as a room or a corridor, the edge represented by a transition region such as a doorway that links the two semantic constructs. Each semantic node embeds within it a topological graph that constitutes the map at the middle level. The topological graph is a set of nodes, each node representing an image of the higher semantic construct. At the low level the topological graph embeds metric values and relations, where each node embeds the pose of the robot from which the image was taken and any two nodes in the graph are related by a transformation consisting of a rotation and translation. The exploration algorithm explores a semantic construct completely before moving or branching onto a new construct. Within each semantic construct it uses a local feature based exploration algorithm that uses a combination of local and global decisions to decide the next best place to move. During the process of exploring a semantic construct it identifies transition regions that serve as gateways to move from that construct to another. The exploration is deemed complete when all transition regions are marked visited. Loop detection happens at transition regions and graph relaxation techniques are used to close loops when detected to obtain a consistent metric embedding of the robot poses. Semantic constructs are labeled using a visual bag of words (VBOW) representation with a probabilistic SVM classifier.

I. INTRODUCTION

Mobile robot exploration is a vital cog in the automation of the mapping process. In recent years, lot of work has been done on image based navigation along the lines of appearance based mapping [3] and topological SLAM[4]. Image based navigation algorithms such as [2] have shown a framework for navigating from one node to another in a topological map based based on images. However these methods do not describe in detail the process of automating the map construction.

While range sensor based exploration has been well understood, the amount of literature on vision based exploration is indeed sparse. One of the first papers in this area [1] adapted the frontier exploration technique [5] to occupancy grid maps constructed through a stereo camera. The maps created had a metric representation. Later frontier exploration technique dovetailed to an image based mapping technique

in [2]. Very recently an image based exploration method that reduced the number of nodes in the topological graph through a combination of local and global decision making strategies was presented in [14] by current authors. Apart from these there isn't any other literature on exploration that is purely guided by visual perception.

The above vision based exploration approaches utilized the image as a provider of dense range information as in [1] or as a characterization of the vicinity around the robot pose as in [2] [14]. It is possible to glean lot more from an image, for example one can obtain the semantic construct in which the robot operates and use this higher level understanding to formulate an effective exploration strategy. Such semantic understanding of places is particularly apt in an indoor and personal robotic context where the robot can communicate with humans through such semantic constructs, take commands from humans in terms of such constructs and on the whole facilitate better interaction between robots and humans.

In this paper we come up with a strategy that provides the robot with hybrid understanding of its surroundings from the lower metric characterizations to higher semantic recognition. The robot explores and constructs a hybrid map that reflects such an understanding. At the highest or top most level the map is a graph whose nodes are semantic constructs such as labs or corridors and the edges represent transition regions (TRs) such as doorways or intersections. At the intermediate or middle level each semantic construct is further detailed as a topological graph of images. At the lowest level the topological graph embeds metric values and relations, where each node embeds the pose of the robot from which the image was taken and any two nodes in the graph are related by a transformation consisting of a rotation and translation. The exploration algorithm first explores the current semantic construct completely before moving to another. Within a semantic construct it builds a topological graph of images by adapting our earlier method delineated in [14]. During this process it identifies gaps as possible gateways of moving from one semantic construct to another. Out of these gaps only few are valid gateways or TRs such as doorways. Gaps are identified using laser data and the confirmation of a gap as a valid TR is through image data. Once a semantic construct is considered completely explored the algorithm moves through one of the TRs to begin exploration of the new construct. The exploration terminates when all gaps are visited and graph relaxation techniques are used to close loops when detected to obtain a consistent metric embedding of the robot poses. Semantic labeling or classification of an

This work is supported by grants from MCIT, India

The authors thank Sathish and Gururaj for their assistance in working on P3DX and their valuable suggestions

Aravindhan K Krishnan is a Graduate Student and Madhava Krishna is a faculty at Robotics Research Center, IIIT Hyderabad
mkrishna@iiit.ac.in

image is through a probabilistic SVM classifier that runs over a bag of words characterization of the image.

The advantages of the method are as follows. Since TRs represent gateways to move from one semantic place to another, loop detection and closure can be done solely based on these TRs, instead of comparing a currently acquired image with all previously obtained images. This significantly reduces the number of comparisons and the computations thereof as well as the number of false loop detections due to reduction in amount of image data being used for such task. Secondly and perhaps most importantly, understanding the larger semantic context within which the robot operates can result in context specific exploration. For example the exploration strategy for a room and corridor can differ if the higher semantic understanding is present.

The novelties of this work arise from advantages mentioned in the previous paragraph, and the advantages that accrue from a semantic understanding in home robotics setting mentioned earlier. Also, this is one of the first efforts that demonstrates a vision-based exploration algorithm that builds hybrid maps.

II. RELATED LITERATURE

Range sensor based exploration became popular in the last decade due to the frontier approach [5] and was further extended to a multi-robotic framework in [6], [7]. [15] presented a decision theoretic approach to multi robot exploration that also included localizing one robot in the map built by the other as an aid to the exploration process. Later Sawhney and others [8] came up with a new per-time visibility metric using which they could explore an unknown area in faster time.

With vision as the primary sensing modality the work of Sim and Little [1] was first of its kind that circumvented the requirement of range sensors for exploration. Frontiers were computed from occupancy grid map just as in [5] and the best frontier to move to was decided through a cost function that trades off distance to reach a frontier with the information gained at that place. Later an image based exploration approach was presented in [2] that once again computed frontiers from images. More specifically they computed frontiers as horizons which are detected as the end of the ground floor segmented from the image. Since this method relies on the assumption that the ground is flat and of a similar texture when compared with other obstacles around its reliability where the ground is undulating such as in outdoors could be in question. Very recently an image based exploration method that reduced the number of nodes in the topological graph through a combination of local and global decision making strategies was presented in [14].

In recent years there have been a lot of approaches tackling the problem of visual topological SLAM [4], visual homing [9] and loop detection [10] where the robot is either guided to acquire the images for learning as in [9] or made to move along predefined paths [10]. An exploration algorithm in such scenarios could limit user intervention as well as extend the range of robot operation as it dynamically expands its

workspace of operation acquiring more images for learning as well as for mapping.

III. METHODOLOGY

This work extends on the topological exploration algorithm of [14] and builds a semantic understanding of the environment. The topological exploration proceeds by a combination of local and global decision making and explores a particular semantic construct, thereby creating a topological graph of images. Further explanation and results for algorithm can be found in [14]. In this paper, we explain how the topological exploration algorithm is adapted within the context of semantic exploration. We then delineate the strategy for identifying gaps and go on to briefly explain how the true gaps are identified as TRs through the probabilistic SVM classifier that learns class labels from images that are characterized through the visual bag of words (VBOW) paradigm [11]. We then describe how loops are detected from TRs and closed by graph relaxation.

A. Semantic Exploration for Semantic Mapping

The semantic exploration strategy is built over the topological exploration strategy as explained below. Each image acquired at a node is classified in terms of previously learned class labels, namely labs (LAB), corridors (CORR), transition regions (TRs), hard to classify (HTC). HTC was included because, when a robot is moving autonomously, not every image will have a distinctive or characteristic view of the place. This was highlighted in [16]. During the exploration process, gaps are identified through the laser readings in a manner explained in the section III-B below. The gaps are represented by their midpoints and identified by their (x,y) co-ordinates and the node in which they are seen at the closest range. Once the current semantic class/construct is sufficiently explored as seen by saturation of the weights [14], each of the gap is visited and checked if it is a TR such as a doorway. Visiting a gap becomes trivial as the shortest path from the current node and the node where the gap was seen in the topological graph can be found. Upon reaching the gap, an image is taken and VBOW+SVM is used to classify the image. Those images (correspondingly gaps) that are not classified as TRs are rejected as false gaps. The robot then picks up one of these TRs to explore, moves across this TR, identifies the new semantic construct and begins exploration of this construct through its underlying topological exploration algorithm. The process continues till all valid TRs are marked visited at which time the exploration halts.

If the semantic construct is a corridor the underlying topological exploration could have the weight of all nodes saturated before the corridor is fully explored due to strong perceptual aliasing due to lack of differentiating features along a corridor. Hence when the robot reaches one end of the corridor, there is no next best node (global decision making as explained in [14]) to choose because of saturation of weights and the exploration stops. To overcome this, knowledge of the higher semantic understanding of the

environment is made use of. Once an image is classified as a corridor the exploration strategy is to simply move towards both the ends of the corridor. The rest of the procedure consists of identifying all other gaps that occur during the exploration and marking out the actual TRs amongst those.

B. Detection of Transition regions (TRs)

TR detection occurs in two steps - Gap identification and Gap verification. For gap identification, the laser readings (r, θ) are converted to (x, y) and the points are clustered (figure 1(a)) based on euclidean distance and line segments are fit for each cluster. Line segments with similar slope are then grouped. Figure 1(b) shows line segments of the same group marked with the same colour. The connecting line segment between adjacent pair of line segments within a group are considered as potential TRs and a visibility check is made. Visibility check is done by generating points on the connecting line segment and computing the corresponding r, θ for those points. This r , is compared with the actual laser reading r' at θ . $(r' - r) > threshold$ for all points implies that the gap is visible and hence could be a transition region.

Figure 1(b) shows an example where the visibility check helps in discarding false gaps. The points generated on the connecting line segments are shown in black and the final gap detected. This approach sometimes results in false gaps, for which an example is shown in figure 1(c). The line segments constituting the gap is superimposed on the image in red. The connecting line segment corresponding to false gap is shown in blue. To eliminate these false gaps, we do a gap verification step.

Gap verification occurs after a room has been explored completely. The robot visits the gaps identified(in the first step) within a room at a close proximity and takes an image at that position to verify if it is a TR or not. Figure 1(c) is the image taken at one of the false gaps. Figure 1(d) is the image taken at a transition region. The probabilistic SVM classified 1(c) as a lab ($\Pr(\text{LAB})= 0.43$, $\Pr(\text{TR}) = 0.21$) and 1(d) as a TR ($\Pr(\text{LAB}) = 0.23$, $\Pr(\text{TR}) = 0.46$). Thus false gaps are completely eliminated by the second step of our approach. Once the gap has been verified, the robot moves through the gap(TR), turns 180° and takes an image from the other side of the TR. The image of the TR thus taken and the midpoint of the TR are used to detect loops. Loop detection is explained in detail in Section III-D. The gap identification method used here is geometric and works well for structured indoor environments. To detect gaps in complex and unstructured environments could entail some form of machine learning techniques as in [12].

C. Semantic Classification of Images

SURF feature descriptors in the training set are extracted and a dictionary of words [11] comprising these features is formed. The frequency of occurrence of these words in every image of the dataset is computed and a vector of such frequencies is formed. A probabilistic SVM uses this vector of frequencies as the lower dimensional input representation of an image and the associated class label as the output vector

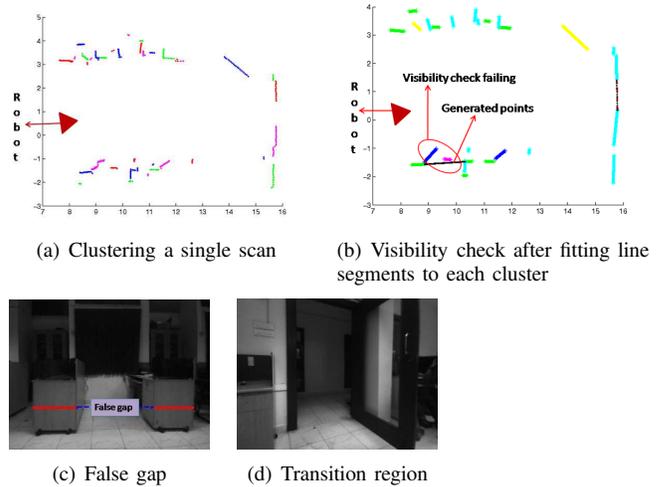


Fig. 1. Detection of Transition Regions

TABLE I
VBOW+SVM PERFORMANCE

Class Label	No. of train images	No. of test images	Accuracy %
Lab	900	160	100
Corridor	683	123	100

and trains over this input-output pair to form class boundaries and obtain the probability of each image belonging to a semantic class. Upon presenting a query image obtained during exploration the trained SVM outputs the probabilities of each semantic class for this query image and the class with the highest probability is the semantic class or construct of that image. The exploration algorithm makes use of this probabilistic SVM and bag of words combine to classify an image online in terms of previously learned semantic class labels.

Transition probabilities along the lines of [12] could be used to accelerate the identification process when the robot has transitioned to a new construct through a transition region. However in our experiments we did not find any tangible advantage by incorporating transition probabilities and thus the identification proceeds by the probabilities as assigned by the SVM classifier.

The training set comprises of 1583 images collected in a particular floor of our college. To validate the efficiency of VBOW+SVM classifier, we tested it against the images taken from [18]. [18] didn't have the classes TR and HTC. So, we tested against only two classes corridor and lab. We selected 283 images randomly from [18] corresponding to the two classes and achieved 100% accuracy. This proves the efficiency of VBOW+SVM in semantic classification. Results are shown in Table I.

D. Loop Detection and Graph Relaxation

Loop detection and closure occurs at the semantic level through the TRs. As mentioned in section III-B each TR is represented by its midpoint and the uncertainty of the robot

TABLE II
DETECTION OF TRANSITION REGIONS(TRs)

Semantic construct	No. of runs	Total no. of TRs	TRs detected
Lab-1	10	2	2
Lab-2	10	2	2
Corridor	10	3	3
Lab-3	5	1	1

projected onto this point along with the measurement error in form of the innovation covariance matrix S . Whenever a TR is seen during the course of exploration, the NIS distance between the current TR and every other TR seen already is computed. If the NIS distance to the closest mapped TR is within a range ($gate1 < nis < gate2$), then the image corresponding to that particular TR (which was taken by moving through the gap and taking a 180° turn) and the current image are compared for similarity. Similarity between images $I1$ and $I2$ is defined as

$$S(I1, I2) = \frac{\text{No. of matching SURF descriptors}}{\min(\text{SURFcount}(I1), \text{SURFcount}(I2))} \quad (1)$$

where $\text{SURFcount}(I)$ is the number of SURF features in image I . If the similarity check is positive (35%), it implies detection of a loop and graph relaxation is run to distribute the error. Graph relaxation is done by treating the midpoints of TRs as nodes in the graph. A TR seen at two nodes (representing robot poses) in the graph during loop detection, will have edges connecting it to both the nodes, thus forming a loop. Thus running graph relaxation now changes the (x,y) co-ordinates of the robot as well as the TRs. We use the graph relaxation algorithm proposed by [13]. [13] corrects only displacement error and not the orientation error. Our orientation estimates were quite accurate. This is evident from the laser plots of figure 3, where the laser scans corresponding to the TR involving a loop appears parallel. Thus, [13] was good enough for graph relaxation in our setup. [17] discusses a method to correct both displacement and orientation errors, which can be used if the orientation estimates are not accurate

E. Localization

Localization is done by searching for the nearest image (say image 1) in the topological graph. The nearest image is found by the count of the matching SURF descriptors. The image in the adjacent node (say image 2) is taken and the matching descriptors in image 1 and image 2 are triangulated and their world co-ordinates are found upto scale. The descriptors corresponding to the world co-ordinates are matched in the current view and thus we obtain a relationship between 2D images points and 3D world points, which can now be used to find the extrinsic parameters (R and T). This R and T localizes the robot corresponding to the nearest node in the topological graph.

IV. EXPERIMENTAL RESULTS

All experiments were carried out on a P3DX robot with a wide angle stereo and SICK laser in an environment spanning labs and corridors. We used a wide angle stereo (for far feature strategy in topological exploration) and so it was sufficient for us to take only 3 images at a node. We have presented results showing exploration, loop detection and closure. Semantic grouping of nodes to form a semantic map is also shown. Obstacle avoidance was done by enabling VFH in Player library.

A. An Exploration run

Figure 2 shows the sequence of images captured during exploration, and the probabilities of their corresponding class labels is also shown. Figure 3(a) shows the path taken by the robot during exploration. Initially the robot is in LAB-1 and does a far-feature based topological exploration, and finds TRs. Figures 2(a) and 2(b) shows the images captured in the lab along with their probabilities. 3 gaps were detected during the process, which are shown in figures 2(c) and 2(e). These were then given to the gap verification routine which classified 2(c) and 2(d) as TRs. Within a semantic construct there is no confusion between TRs because the odometry error doesn't grow so much that the uncertainty ellipses overlap. Even when it overlaps, a similarity check of the images captured in the TRs helps us in resolving the ambiguity.

The robot then moves out of the two transition regions (TR1 and TR2), turns 180° and takes images. These are the images which will be used for loop detection when the same transition regions are visited from the other side of the transition region.

The robot then decides to further explore by moving out of transition region TR2 (figure 2(d)), marking it as visited. The new semantic construct is a corridor and the robot explores this semantic construct using a "corridor specific exploration strategy". The images taken during exploration of the corridor is shown from figures 2(f) to 2(g). Gap identification in the corridor is quite robust in our approach, because the corridor environment is usually uncluttered, and there are no false gaps. So, we don't do gap verification if the semantic construct being explored is a corridor. The robot identifies two transition regions (TR2 and TR3) within the corridor and decides to move into TR3 as TR2 is already marked as visited. After exploring the corridor, it moves into the unvisited TR3 and starts exploring LAB-2 based on far-feature based strategy. The robot first moves to the right, hits a dead end and finds the next best node (global decision making) and starts exploring towards the other side of the lab. The images taken in the LAB-2 during exploration is shown from figures 2(h) and 2(i). During the course of exploration, it detects a gap whose uncertainty ellipse overlaps with a transition region (TR1) seen already. The current image (figure 2(i)) and the image captured on the other side of the gap (figure 2(j)) are compared which look similar, hence the loop is detected (explained in more detail in section III-D). A transition region is marked visited when

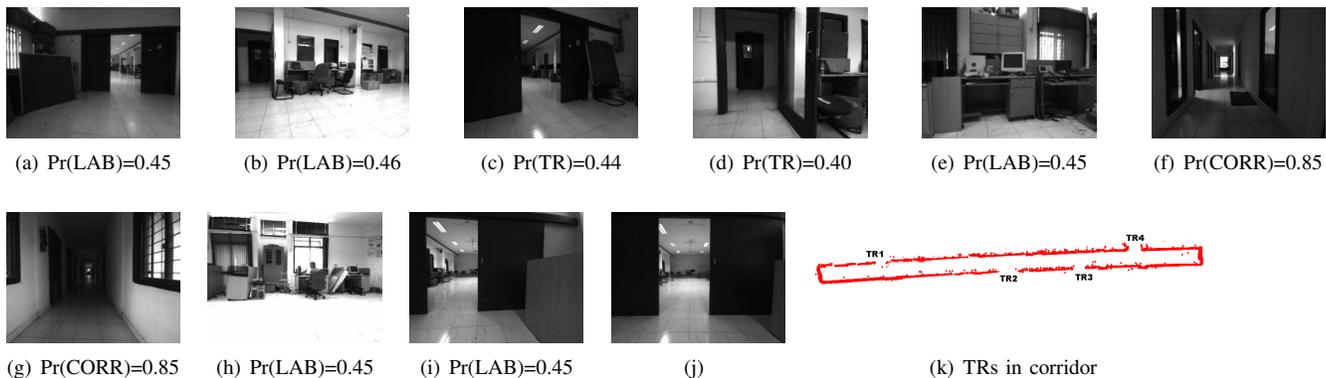


Fig. 2. Exploration

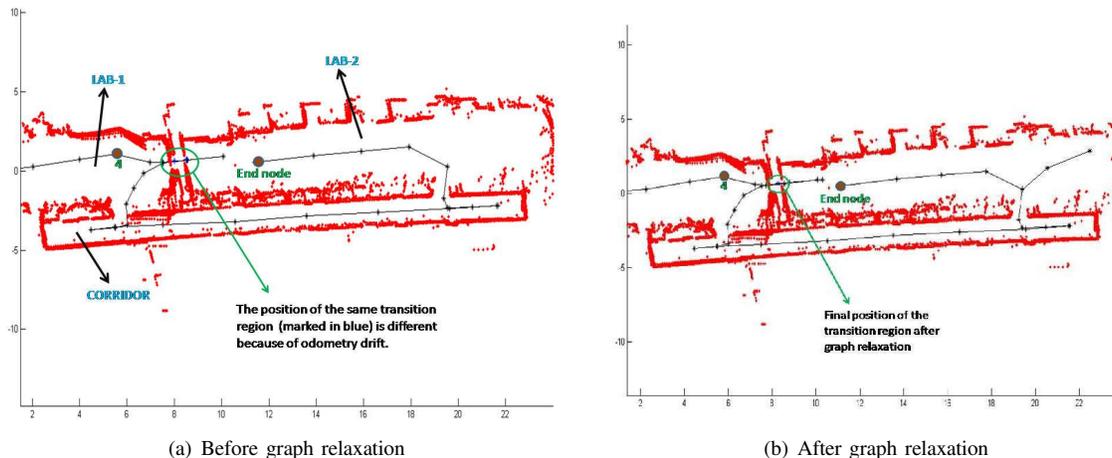


Fig. 3. Loop closure during exploration

a loop is detected on it, thus making TR1 visited. Graph relaxation algorithm is run at this stage and the loop is closed. At this point all transition regions (TR1, TR2 and TR3) have been marked visited and the robot also reaches the end of LAB-2. Hence exploration is terminated.

B. The Role of Transition Regions

Transition Regions are used for terminating the exploration, serve as gateways to move from one semantic construct to another to expand exploration and as well as for loop detection.

The ability to differentiate between TRs within and across semantic constructs without confusion helps us in loop detection at the level of TRs. This is definitely faster than the loop detection methods that compare a current image with most of the previously acquired images during exploration. Here we present arguments as well as empirical results to substantiate about TRs being effective loop detection agents. For being such effective agents they must be detected when present and not be falsely associated with another TR.

Within a semantic construct that is being explored TRs never get falsely associated. In a room construct the number of TRs are very small (one or two at most) and far apart to get easily discerned by both distance and image features. Within a corridor despite being visually similar and being

close enough they are not falsely associated since one can exploit the semantic understanding of being within a corridor. That the TRs occur only on either side of a corridor results in a simple but robust enumeration of such regions. For example in figure 2(k) during the upward journey of a corridor two TRs were detected on the left followed by two on the right. Hence it is evident on the return the same TRs would be detected on the reverse order and in the opposite directions. This simple reasoning scheme effectively exploits the semantic understanding of the corridor to correctly associate TRs even if they are close enough to cause confusions at the level of odometry.

Table II shows detection performance of TRs within a semantic construct. The experiments were conducted by running the exploration algorithm within a particular semantic construct by starting the robot from different places. The table shows experiments in 3 labs and a corridor, the number of transition regions involved in each of the semantic construct and the number of TRs detected. It can be seen from the table that all the TRs within a semantic construct were detected.

The question then arises how well are the TRs detected across constructs. This pertains to loop detection situation where a TR seen previously from one semantic construct is now viewed from a different construct. We explain the

loop detection process and present empirical results in the subsequent subsection.

C. Loop closure at transition regions

Figure 3 shows loop detection and closure occurring in the experiment described in the above section. Figure 3(a) shows loop detection. The same TR (marked in blue) is seen at node 4 and end node. The (x,y) co-ordinate of the TR seen at both nodes are not same because of odometry errors. The odometry drift can be seen in the laser plot in figures 3(a) and 3(b). The NIS distance between the TRs seen at node 4 and end node is within a range, hence invoking the image comparison as described in section III-D. The images used for comparison are shown in 2(j) and 2(i). The images match and hence a graph relaxation algorithm is run. We modelled the uncertainty of nodes/TRs as an ellipse and not as a circle as discussed in [13]. This extension was trivial. Figure 3(b) shows the corrected graph after running the graph relaxation algorithm. Graph relaxation was run after detecting loop.

Few runs were taken by manually guiding the robot to test our loop detection method. In all these runs we found the loops were effectively detected through TRs thus showing their ability as effective loop closure agents. The corridor was blocked on two ends by us in all our experiments because of height discontinuities. This is why both the ends of the corridor appear as dead ends in our results. Height discontinuities in our corridors prevents our testing in larger environments, which would be taken up in future. Nonetheless the framework for semantic exploration presented in this paper with vision as the chief sensing modality is novel and appropriate for indoor and home robotic settings.

D. Deriving a semantic map from a topological map

Once the robot crosses a TR and enters a new semantic construct, all images/nodes seen in the previous semantic construct are labelled with the same semantic label and hence a semantic map can be derived by grouping nodes of semantic constructs and their connecting TRs.

The semantic map for the exploration explained in IV-A is shown in figure 4(b).

V. CONCLUSIONS

An novel exploration framework was presented to build a hybrid map of the environment - topological map of images at a lower level and semantic grouping of nodes in the topological graph to form a semantic map at a higher level. This was done using vision as the main sensing modality with laser being minimally used for obstacle avoidance and gap identification. Loop detection was done at the level of transition regions between two semantic constructs and the loops were closed using graph relaxation techniques thus averaging out the odometry and measurement errors. Results show the efficacy of the proposed exploration framework. The scale of the test environment is bigger than a typical home setting which suggests that the framework is ideal for home robotics. The current framework is applicable to larger environments as well, perhaps with more semantic

categories like junctions/intersections etc. Our future work would address this.

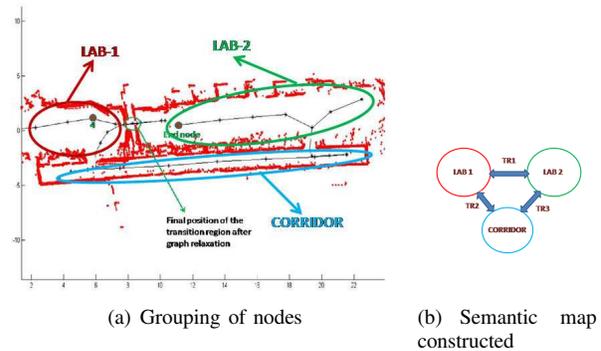


Fig. 4. Semantic Mapping

REFERENCES

- [1] R. Sim and J. J. Little. Autonomous vision-based exploration and mapping using hybrid maps and Rao-blackwellised particle filters. In *IOS*, pages 2082–2089, 2006.
- [2] D. Santosh, S. Achar, and C. V. Jawahar. Autonomous image-based exploration for mobile robot navigation. In *ICRA*, pages 2717–2722, 2008.
- [3] M. J. Cummins and P. M. Newman. Fab-map: Probabilistic localization and mapping in the space of appearance. *International Journal of Robotics Research*, 27(6):647–665, 2008.
- [4] A. Angeli, S. Doncieux, J.-A. Meyer, and D. Filliat. Incremental vision-based topological slam. In *IOS*, pages 1031–1036, 2008.
- [5] Y.B Yamauchi. A frontier-based approach for autonomous exploration. In *CIRA*, page 146, 1997.
- [6] W. Burgard, M. Mooors, C. Stachniss and F. Schneider. Coordinated multi-robot exploration. In *IEEE Transactions on Robotics* pages 376–378, 2000
- [7] R. Simmons, D. Apfelbaum, W. Burgard, D. Fox, S. Moors, S. Thrun. Coordination for multirobot exploration and mapping. in *Proc. of the National Conf. on Artificial Intelligence(AAAI)*, 2000
- [8] R. Sawhney, K.M. Krishna and K. Srinathan. On fast exploration on 2d and 3d terrains with multiple robots. In *AAMAS* pages 73-80, 2009
- [9] D. Filliat. Interactive learning of visual topological navigation. In *IOS* pages 248-254, 2008
- [10] A. Angeli, S. Doncieux, J.A. Meyer, D. Filliat. Real time visual loop-closure detection. In *ICRA* pages 1842-1847, 2008
- [11] G. CSurka, C. Dance, L. Fan, J. Willamowski Visual categorization with bag of keypoints. In *Workshop on Statistical Learning in Computer Vision* 2004
- [12] Alex Rottmann, Oscar Martinez Mozos, Cyrill Stachniss, Wolfram Burgard. Semantic Place classification of Indoor Environments with Mobile Robots using Boosting. In *AAAI* 2005
- [13] T. Duckett, S. Marsland, J. Shapiro. Learning globally consistent maps by relaxation. In *ICRA* 2000
- [14] Aravindhan K Krishnan, K Madhava Krishna, Supreeth Achar. Image based Exploration in Indoor Environments using Local Features. In *AAMAS* 2010.
- [15] J. Ko, B. Stewart, D. Fox, K. Konolidge, B. Limketkai A practical decision theoretic approach to multi-robot mapping and exploration In *IOS* pages 415-421, 2003
- [16] Jianxin Wu, Henrik I. Christensen, James M. Rehg Visual Place Categorization: Problem, Dataset, and Algorithm In *IOS* 2009
- [17] U. Frese, P. Larrson, T. Duckett A multilevel relaxation algorithm for simultaneous localization and mapping In *IEEE Transactions on Robotics* 2005
- [18] A. Pronobis and B. Caputo COLD: Cosy Localization Database In *International Journal of Robotics Research* May 2009