



Queensland University of Technology
Brisbane Australia

This may be the author's version of a work that was submitted/accepted for publication in the following source:

Lowry, Stephanie, Suenderhauf, Niko, Newman, Paul, Leonard, John, Cox, David, Corke, Peter, & Milford, Michael
(2016)

Visual place recognition: A survey.

IEEE Transactions on Robotics, 32(1), pp. 1-19.

This file was downloaded from: <https://eprints.qut.edu.au/222264/>

© Consult author(s) regarding copyright matters

This work is covered by copyright. Unless the document is being made available under a Creative Commons Licence, you must assume that re-use is limited to personal use and that permission from the copyright owner must be obtained for all other uses. If the document is available under a Creative Commons License (or other specified license) then refer to the Licence for details of permitted re-use. It is a condition of access that users recognise and abide by the legal requirements associated with these rights. If you believe that this work infringes copyright please provide details by email to qut.copyright@qut.edu.au

Notice: *Please note that this document may not be the Version of Record (i.e. published version) of the work. Author manuscript versions (as Submitted for peer review or as Accepted for publication after peer review) can be identified by an absence of publisher branding and/or typeset appearance. If there is any doubt, please refer to the published source.*

<https://doi.org/10.1109/TRO.2015.2496823>

Visual Place Recognition: A Survey

Stephanie Lowry, Niko Sünderhauf, Paul Newman, *Fellow, IEEE*, John J. Leonard, *Fellow, IEEE*,
David Cox, Peter Corke, *Fellow, IEEE*, and Michael J. Milford, *Member, IEEE*

Abstract— Visual place recognition is a challenging problem due to the vast range of ways in which the appearance of real-world places can vary. In recent years improvements in visual sensing capabilities, an ever-increasing focus on long-term mobile robot autonomy, and the ability to draw on state of the art research in other disciplines – particularly recognition in computer vision and animal navigation in neuroscience – have all contributed to significant advances in visual place recognition systems. This paper presents a survey of the visual place recognition research landscape. We start by introducing the concepts behind place recognition – the role of place recognition in the animal kingdom, how a “place” is defined in a robotics context, and the major components of a place recognition system. We then survey visual place recognition solutions for environments where appearance change is assumed to be negligible. Long term robot operations have revealed that environments continually change; consequently we survey place recognition solutions that implicitly or explicitly account for appearance change within the environment. Finally we close with a discussion of the future of visual place recognition, in particular with respect to the rapid advances being made in the related fields of deep learning, semantic scene understanding and video description.

Index Terms—Visual Place Recognition.

I. INTRODUCTION

VISUAL place recognition is a well-defined but extremely challenging problem to solve in the general sense; given an image of a place, can a human, animal or robot decide whether or not this image is of a place it has already seen? Whether referring to humans, animals, computers or robots, there are some fundamental things a place recognition system must have and must do. Firstly, a place recognition system must have an internal representation – a map – of the environment to compare to the incoming visual data. Secondly, the place recognition must report a belief about

whether or not the current visual information is from a place already included in the map, and if so, which one. Performing visual place recognition can be difficult due to a range of challenges; the appearance of a place can change drastically (see Fig. 1), multiple places in an environment may look very similar, a problem known as perceptual aliasing, and places may not always be revisited from the same viewpoint and position as before.

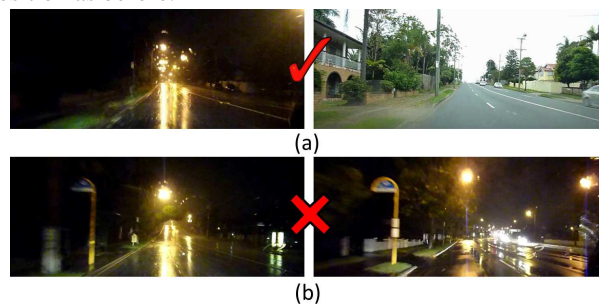


Fig. 1. A visual place recognition system must be able to (a) successfully match very perceptually different images while (b) also rejecting incorrect matches between aliased image pairs of different places.

In robotics, this research topic is highly relevant given the ever increasing focus on long term mobile robot autonomy and rapid improvements in visual sensing capabilities and cost. Vision is the primary sensor for many localization and place recognition algorithms [1]–[19]. Place recognition is also a growing research field, as evidenced by citation analyses and a number of dedicated place recognition workshops at recent and upcoming robotics and computer vision conferences including the *International Conference on Robotics and Automation* (2014, 2015) and the *IEEE Conference on Computer Vision and Pattern Recognition* (2015). The problem of persistent place recognition has also formed a regular component of many more general workshops including the long-running ICRA workshop on Long-Term Autonomy (2011 – 2014).

Our aim in writing this survey article is to provide a comprehensive review of the current state of place recognition research that is relevant both to robotics and other fields of research including computer vision and neuroscience. The timing for such a survey is particularly fortuitous given major events across these related fields: for example, the almost universal usage of deep learning techniques in state of the art recognition systems in computer vision, and the 2014 Nobel Prize in Physiology or Medicine award to Edvard Moser, May-Britt Moser and John O’Keefe, who discovered the key representations of place in the mammalian brain. This paper provides an overview of the place recognition problem and its

Submitted for review on 18 March 2015. This research was supported by Microsoft and the Australian Centre for Robotic Vision.

S. Lowry, N. Sünderhauf, P. Corke and M. J. Milford are with the Australian Centre of Robotic Vision, School of Electrical Engineering and Computer Science, Queensland University of Technology, QLD 4000 Australia (e-mail: firstname.lastname@roboticvision.org).

P. Newman is with the Mobile Robotics Group, Department of Engineering Science, University of Oxford, OX1 3PJ UK (e-mail: pnewman@robots.ox.ac.uk).

J. J. Leonard is with the Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, MA 02139 USA (e-mail: jleonard@csail.mit.edu).

D. Cox is with the Department of Molecular and Cellular Biology, the School of Engineering and Applied Science, and the Center for Brain Science, Harvard University, MA 02138 USA (e-mail: davidcox@fas.harvard.edu).

relationship with many major robotics research fields including SLAM, localization, mapping, and recognition. Because of the increasing focus in the research community on long term robot autonomy in challenging environments, we also provide a particular focus on the problem of lifelong visual place recognition for robots.

II. THE CONCEPT OF PLACE IN ROBOTICS AND THE NATURAL KINGDOM

The problem of navigation and place recognition has a venerable tradition in psychology and neuroscience. In 1948, the research of Tolman [20] on rats navigating mazes motivated him to propose the cognitive map – a mental representation of the world with information about relationships between places that animals gradually learn. The concept of the cognitive map, while not without its critics [21], [22], has been influential not only in psychology and neuroscience, but also areas such as urban planning, where Lynch [23] proposed that the elements of a cognitive map be paths, edges, nodes, districts and landmarks, and in robotics, where mapping approaches have been inspired by the cognitive map [24], [25], and by its successor, the spatial semantic hierarchy [26].

With the development of techniques to record neural activity in the brain of animals [27] came the identification of place cells in the rat hippocampus by O’Keefe and Dostrovsky [28]. Place cells fire when the rat is in a particular place in the environment (see Fig. 2(a)), and the population of place cells cover the entire environment [29], [30]. Furthermore, if a rat moves from one environment to another, the same place cells can be used to represent multiple different environments. O’Keefe and Conway [31] proposed that these place cells form a part of Tolman’s cognitive map. The understanding about the relationships between neural activity and places in the world was extended by the discovery of head direction cells in the dorsal presubiculum [32] and of grid cells [33] in the medial entorhinal cortex (MEC). Head direction cells fire when an animal turns its head in a particular direction relative to its body, while grid cells fire in multiple places in the environment, in such a format that their firing fields form a regular grid (see Fig. 2(b)).

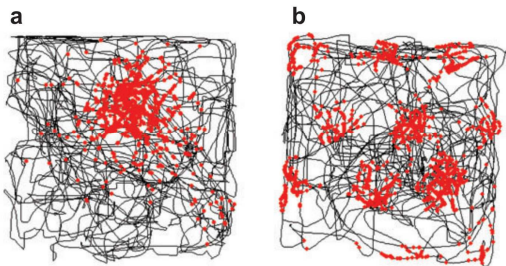


Fig. 2. Neuroscience experiments have shown that the brains of animals such as rats contain place cell and grid cell neurons. Each place cell fires strongly at one location in an environment, while each grid cell fires at multiple, regularly spaced locations. This figure shows the firing locations of (a) a place cell and (b) a grid cell placed over the path of an animal in a square environment (from [34]).

Place recognition, as observed via the firing of place cells, is triggered by both sensory cues and self-motion [29]. Studies

with rats show that place cell firing is initially based on self-motion, but if the environment is changed - by altering the distance between start and end goals, for example - the place cell will update to the correct location according to the external visual landmarks [35], [36]. The correction may occur smoothly or abruptly, depending on the size of the mismatch.

Many of the same concepts arise in robotics. Most robots have access to external observation data as well as self-motion information. Topological and metric relationships between places are used in combination with sensory cues to determine the most likely place, similar to the neuronal firing of the place cells. Fig. 3 presents a schematic of a visual place recognition system. Visual place recognition systems contain three key components – an *image processing* module to interpret the incoming visual data; a *map* that maintains a representation of the robot’s knowledge of the world; and a *belief generation* module, which uses the incoming sensor data in combination with the map to make a decision about whether the robot is in a familiar or novel place. A place recognition system may also use motion or transition information to inform the belief generation process. Furthermore, most place recognition systems are designed to operate online, and thus must update the map accordingly.

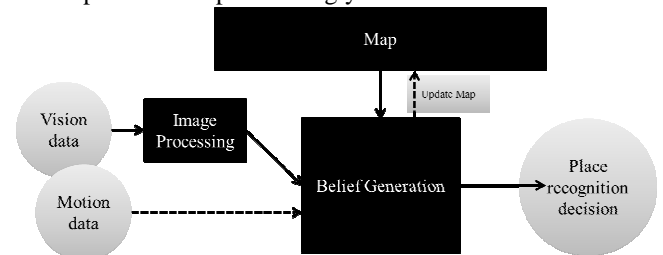


Fig. 3. Schematic of a visual place recognition system. Incoming visual data is processed by the *image processing* module. The robot’s knowledge of the world is stored in the *map*. The *belief generation* module decides whether the current visual data matches a previously stored place. Motion information is also often included, and the map may be continually updated during operation.

This paper discusses what qualifies as a place in the context of robotic navigation. It then looks at the three key modules that make up the place recognition system: the image processing module, the mapping framework, and the belief generation module. The paper then turns to the problem of changing environments. It revisits each of the modules – the image processing module, the mapping module and the belief generation module – and investigates how each has to be adapted to incorporate the notion of appearance change into the place recognition system’s model of the world.

III. WHAT IS A PLACE?

The concept of places in robotics is motivated by the challenges of robotic navigation and mapping. A real robot has fallible sensors and actuators and it is challenging to build a metrically accurate map of the world, and to maintain self-localization within such a representation. The combination of both these goals, known as Simultaneous Localization and Mapping (SLAM) [37]–[41], is even more difficult to consistently achieve.

An alternative approach is to use a “relational map, which is

rubbery and stretchy, rather than to try to place observations in a 2-D coordinate system” (Brooks, [40]). Such a topological map is conceptually similar to the biological notion of a cognitive map, and uses nodes to represent the possible places in the world and edges to represent the possible paths between these places. Robot navigation is reduced to following these edges between nodes and the places represent key intersections or decision points between routes [42], [43] as well as desirable end goals.

This topological approach to navigation is not without difficulties. The robot has to associate these abstract routes and places with physical places and paths, and the complex relationship between the robot sensors, the robot controls, and the robot’s topological and metric interpretations of the world need to be defined [26]. Another issue is how a robot can generate topological maps. If the robot has access to a metric gridmap of the environment, it can extract topological information, emphasizing relevant navigation information like open spaces and passageways [44]. Alternatively, a topological map can be created by a robot from visual and transition information.

The definition of a place depends on the navigation context, and may either be considered as a precise position – “a place describes part of the environment as a zero-dimensional point” (Kuipers, [26]), or as a larger area – “a place may also be defined as the abstraction of a region” where a region “represents a two-dimensional subset of the environment” (Kuipers, [26]). A place can be a fairly large two-dimensional physical area – for example, a room in a building might in some cases qualify as a single place, while in other cases it might contain many different places. A region could also be defined as a three-dimensional area, depending on the requirements of the environment or robot. Unlike a robot pose, a place does not have an orientation, and an ongoing challenge in place recognition is pose invariance – ensuring recognition regardless of the orientation of the robot within the place.

The location of each place – whether a one-dimensional point or a larger region – can be selected based on spatial or temporal density. In this approach, a new place is added according to a particular time step, or when the robot has travelled a certain distance. Alternatively, a place can be defined in terms of its appearance. Kuipers and Byun [25] defined a place as somewhere distinctive relative to other nearby locations, according to some associated sensory information known as a place signature or place description. While the distinctiveness criterion is not always required, a topological place is defined as having a certain appearance configuration [45], [46] and the physical bounds of a place occur where the appearance changes significantly, called a “gateway” [47].

This qualitative concept of topological places as regions that are visually homogeneous needs to be quantified – that is, how can a place recognition system actually segment the world into distinct places? Ranganathan [48] noted that there are similarities with the problem of change-point detection in video segmentation [49], [50], and used change-point detection algorithms such as Bayesian surprise [50] and

segmented regression [51] to define places within a topological map [48], [52]. These methods create a new place when the current appearance (determined from the sensor measurements) is unlikely according to the current model of the environment, and therefore a new model is required (see Fig. 4). Similarly, Korrapati, Courbon et al. [53] used Image Sequencing Partitioning (ISP) techniques to group visually similar images together as topological graph nodes, while Chapoulie, Rives et al. [54] combined Kalman filtering with the Neyman-Pearson Lemma. Murphy and Sibley [55] combined dynamic vocabulary building [56] and incremental topic modelling [57] to continually learn new topological places in an environment, and Volkov, Rosman et al. [58] used coresets [59] to segment the environment. Topic modeling, coresets, and Bayesian surprise techniques can also be used for other aspects of robotic navigation, such as summarizing a robot’s past experience [60]–[62], or determining exploration strategies [63].

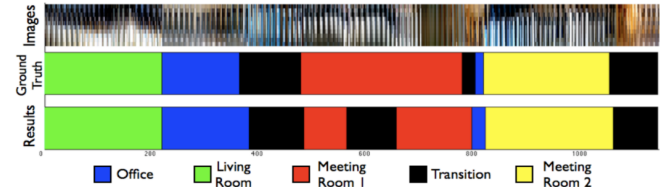


Fig. 4. Topological place recognition systems segment the image stream into places based on the visual information. When a significant change is observed, a new place will be created. In this example (from [48]), the incoming image stream (top row) is segmented based on the detected change points. The detected places (bottom row) match closely to the different rooms shown by the ground truth location (middle row).

Appearance-based and density-based place selection methods are practical to implement as they depend on measurable quantities such as distance, time or sensor values [64]. An ongoing challenge is the enhancement of appearance information with semantic labels such as “door” or “intersection” so places can be selected online based on their value as decision points. The addition of semantic data to maps can improve planning and navigation tasks [65], and requires place recognition to be linked with other recognition and classification tasks, especially scene classification and object recognition. These relationships are symbiotic – place recognition can improve object detection by providing contextual priming for object detection as well as contextual priors for object localization [66], and conversely, object recognition can also aid place recognition [67]–[70], particularly in indoor environments where the function of a place such as “kitchen” or “office” can be inferred from the objects within it, and used to infer the location from a labeled semantic map [71].

IV. DESCRIBING PLACES: THE IMAGE PROCESSING MODULE

Visual place description techniques fall into two broad categories; those that selectively extract parts of the image that are in some way interesting or notable, and those that describe the whole scene, without a selection phase. Examples of the first category are local feature descriptors such as SIFT [72] and SURF [73]. Local feature descriptors first require a

detection phase which determines the parts of the image to retain as local features (see Fig. 5(a)). In contrast, an example of the second category are global or whole-image descriptors such as Gist [74], which do not have a detection phase but process the whole image regardless of its content (see Fig. 5(b)).



Fig. 5. Visual place description techniques fall into two broad categories. (a) Interesting or salient parts of the image are selected for extraction, description and storage. For example, SURF [73] extracts interest points in an image for description. The number of possible features may vary depending on the number of interest points detected in the image. The red circles are interest points selected by SURF within this image. (b) The image is described in a pre-defined way without first detecting interest points. Whole-image descriptors such as Gist [74], [75] divide an image into blocks as shown by the red lines and processes each block regardless of its content.

A. Local feature descriptors

The development of the local feature method Scale-Invariant Feature Transforms or SIFT [72] led to its widespread use in place recognition [76]–[83]. As other local feature detection and description methods were developed, they too were applied to the visual localization and place recognition problem. For example, Ho and Newman [84] use Harris affine regions [85], Murillo, Guerrero et al. [86] and Cummins and Newman [87] use Speeded-Up Robust Features (SURF) [73], while FrameSLAM [2] uses CenSurE [88]. Since local feature extraction consists of two steps – detection followed by description – it is not uncommon to combine different techniques for each. For example, Mei, Sibley et al. [89] use the detection technique FAST [90] to find keypoints in the image, which are then described by SIFT descriptors. Similarly, Churchill and Newman [15] use FAST extraction combined with BRIEF [91] descriptors.

Each image may contain hundreds of local features, and directly matching image features can be inefficient. The bag-of-words model [92], [93] increases efficiency by quantizing local features into a vocabulary that can be compared using text retrieval techniques [94]. The bag-of-words model partitions a feature space, such as SIFT or SURF descriptors, into a finite number of visual words (see Fig. 6). A typical vocabulary contains 5000 – 10,000 words, but a vocabulary as large as 100,000 words has been used for place recognition by FAB-MAP 2.0 [87]. For each image, every feature is assigned to a particular word, ignoring any geometric or spatial structure, thereby allowing images to be reduced to binary strings or histograms of length n , where n is the number of words in the vocabulary.

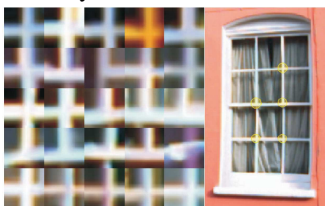


Fig. 6. A bag-of-words model clusters similar features into a single visual word to make recognition more efficient and straightforward. This image (from [6]) shows examples of image patches all corresponding to a single visual word. This word matches window frame crosspieces.

Images described using the bag-of-words model can be efficiently compared using binary string comparison such as a Hamming distance or histogram comparison techniques. Vocabulary trees [95] can make the process for large-scale place recognition even more efficient. Originally proposed for object recognition, vocabulary trees use a hierarchical model to define words, an approach that enables faster lookup of visual words and the use of a larger and thus more discriminating vocabulary. Localization systems that use the bag-of-words approach include [82], [84], [87], [96], [97] and many others.

Because the bag-of-words model ignores the geometric structure of the place it is describing, the resulting place description is pose invariant: the place can be recognized regardless of the position of the robot within the place. However, the addition of geometric information to a place has been shown to improve the robustness of place matching, particularly in changing conditions [14], [87], [98]–[100]. These systems may assume a laser sensor is available for 3D information [98], use stereo vision [14], epipolar constraints [100], [101], or simply define the scene geometry according to the position of the elements within the image [102], [103]. The trade-off between pose invariance – recognizing places regardless of the robot orientation – and condition invariance – recognizing places when the visual appearance changes – has not yet been resolved, and is a current research challenge in place recognition.

The bag-of-words model is typically pre-defined based on features extracted from a training image sequence. This approach can be limiting as the resulting model is environment-dependent and needs to be re-trained if a robot is moved into a new area. Nicosevici and Garcia [56] propose an online method to continuously update the vocabulary based on observations, while still being able to match prior observations with future observations. As a result, a bag-of-words model can be used without requiring a pre-training phase, and can adapt to the environment, out-performing pre-trained models despite requiring less *a priori* knowledge [56].

B. Global descriptors

Global place descriptors used in early localization systems included color histograms [5] and descriptors based on principal component analysis [104]. Lamon, Nourbakhsh et al. [105] used a variety of image features – such as edges [106], corners [107] and color patches – combined into a fingerprint of a location. By ordering these features in a sequence between 0° and 360° , place recognition could be reduced to string-matching. These systems used omnidirectional cameras which allowed rotation-invariant matching at each place.

Global descriptors can be generated from local feature descriptors by pre-defining the keypoints in the image – for example, using a grid-based pattern – and then using the chosen feature description method on the pre-selected keypoints. Badino, Huber et al. [108] used whole-image

descriptors based on SURF features known as WI-SURF to perform localization and BRIEF-Gist [109] used BRIEF features [91] in a similar whole-image fashion.

A popular whole-image descriptor is Gist [74], [75] which has been used for place recognition on a number of occasions [110]–[113]. Gist uses Gabor filters at different orientations and different frequencies to extract information from the image. The results are averaged to generate a compact vector that represents the “gist” of a scene.

C. Describing places using local and global techniques

Local and global descriptors each have different advantages and disadvantages. Local feature descriptors are not restricted to defining a place only in terms of a previous robot pose, but can be recombined to create new places that have not previously been explicitly observed by the robot. For example, Mei, Sibley et al. [114] defined places via co-visibility: the system finds cliques in the landmark co-visibility map which define places even when the landmarks have not simultaneously been seen in a single frame, and can outperform standard image-based place recognition [78]. Lynen, Bosse et al. [115] generated a 2D space of descriptor votes where regions of high vote density represent loop closure candidates.

Local features can also be combined with metric information to allow metric corrections to localization [2], [7], [76]. Global descriptors do not have the same flexibility, and furthermore, whole-image descriptors are more susceptible to change in the robot’s pose than local descriptor methods, as whole-image descriptor comparison methods tend to assume that the camera viewpoint remains similar. This problem can be somewhat ameliorated by the use of circular shifts as in [116] or by combining a bag-of-words approach with a Gist descriptor on segments of the image [17], [110].

While global descriptors are more pose dependent than local feature descriptors, local feature descriptors perform poorly when lighting conditions change [117] and are comprehensively out-performed by global descriptors at performing place recognition in changing conditions [118], [119]. Using global descriptors on image segments rather than whole images may provide a compromise between the two approaches, as sufficiently large image segments exhibit some of the condition invariance of whole images, and sufficiently small image segments exhibit the pose invariance of local features. McManus, Uproft et al. [120] used the global descriptor HOG [121] on image patches to learn condition invariant scene signatures, while Sünderhauf, Shirazi et al. [122] used the Edge Boxes object proposal method [123] combined with a mid-level Convolutional Neural Network (CNN) feature [124] to identify and extract landmarks as illustrated in Fig. 7.



Fig. 7. Object proposal methods such as the Edge Boxes method [123] shown here were developed for object detection but can also be used to identify potential landmarks for place recognition. The colored boxes in the images above show landmarks that have been correctly matched between two viewpoints of a scene (from [122]).

D. Including 3D information in place descriptions

The image processing techniques described above are appearance-based – they “model the data directly in the visual domain (instead of making a geometric model)” (Krose, Vlassis et al., [125]). However, in metric localization systems, the appearance-based models must be extended with metric information. Monocular image data is not a natural source of geometric landmarks – “the essential geometry of the world does not ‘pop out’ of images the same way as it does from laser data” (Neira, Davison et al., [126]). While many systems use data from additional sensors such as lasers [98] or RGB-D cameras [127]–[129], geometric data can also be extracted from conventional cameras to allow metric calculation of the robot pose.

Metric range information can be inferred using stereo cameras [2], [130]–[132]. Monocular cameras can also infer metric information using Structure-from-Motion algorithms [133]. Methods include MonoSLAM [7], PTAM [134], DTAM [135], LSD-SLAM [136] and ORB-SLAM [137]. Metric information can be sparse: that is, range measurements are associated with local features such as image patches as in MonoSLAM [7], SIFT features as in Se, Lowe et al. [76], CenSurE features as in FrameSLAM [2], or ORB features [138] as in ORB-SLAM [137]. In contrast, DTAM stores dense metric information about every pixel, and LSD-SLAM maintains semi-dense depth data on the parts of the image containing structure and information. Dense metric data allows a robot to perform obstacle avoidance and metric planning as well as mapping and localization, so fully autonomous vision-only navigation can be performed [16].

The introduction of novel sensors, such as RGB-D cameras, that provide dense depth information as well as image data has spurred the development of dense mapping techniques [70], [127]–[129], [139], [140]. These sensors can also exploit 3D object information to improve place recognition. SLAM++ [70] stores a database of 3D object models and uses this database to perform object recognition during navigation, and uses these objects as high-level place features. Objects have a number of advantages over low-level place features: they provide rich semantic information, and can reduce memory requirements via semantic compression; that is, storing object labels rather than full object models in the map [70].

V. REMEMBERING PLACES: THE MAPPING MODULE

For a place recognition or navigation task, the system needs to refer to a map – a stored representation of the robot’s knowledge of the world – to which the current observation is compared. The map framework differs depending on what data is available and what type of place recognition is being performed. Table I displays a taxonomy of mapping approaches, which depends on the level of physical abstraction in the map, and whether or not metric information is included in the place description. The most concrete mapping framework listed is the topological-metric or topometric map. Although it is possible to have a globally metric map, such maps are only feasible in small geographical areas, and there are mechanisms for fusing topometric maps into globally metric maps [141]. Thus for the purposes of place recognition any globally metric map can be considered as a one-node topometric map.

TABLE I
MAPPING FRAMEWORKS FOR VISUAL PLACE RECOGNITION

Level of map abstraction	Place description type	Comments
Pure image retrieval	Appearance-based	No position information
Topological	Appearance-based	Includes transition information
Topological-metric	Appearance-based	Includes metric information between but not within places
	Sparse metric information (landmark maps)	SLAM system – includes metric information between and within places
	Dense metric information (occupancy grid maps)	

A. Pure image retrieval

The most abstract form of mapping framework for place recognition only stores appearance information about each place in the environment, with no associated position information. Pure image retrieval assumes that matching is based solely appearance similarity and applies image retrieval techniques from computer vision that are not specific to place-based information [3]. Although valuable information is lost by not including relative position information, there are computationally efficient indexing techniques that can be exploited.

A key concern with place recognition is system scalability – as the robot visits more and more places, storage requirements will increase and search efficiency will decrease. As a result, maps need to be designed to ensure large-scale efficiency. If a bag-of-words model is used to quantize the descriptor space, image retrieval can be accelerated using inverted indices; the image ID numbers are stored against the words that appear in the image, rather than the words being stored against the image IDs. Inverted indices allow much quicker elimination of unlikely images, rather than requiring a linear search of all images in the database.

Schindler, Brown et al. [3] used a hierarchical vocabulary tree [95] to achieve efficient visual place recognition of a city-sized dataset (a 20km traversal with around 100 million

features). This paper showed that place recognition performance improves if only the most informative features from each image are used, where information gain is measured using a conditional entropy calculation. Improved place recognition with a reduced feature set was also observed by Li and Koščeká [142].

FAB-MAP 2.0 [87], [143] also used an inverted index with a bag-of-words model to demonstrate visual place recognition across a 1000 km path. While Schindler, Brown et al. [3] used a voting scheme to match locations, FAB-MAP’s probabilistic model that includes negative observations – words that do *not* appear in the image – as well as positive observations requires simplification before the inverted index approach can be applied.

Place recognition can also be made more efficient by using hierarchical searching at the place level as well as at the vocabulary level. Mohan, Gálvez-López et al. [144] selected the most likely environment using co-occurrent feature matrices. Then place matching is performed using only a subset of the previously seen places, reducing the time required for searching.

B. Topological maps

Pure topological maps contains information about relative positions of places but do not store metric information regarding how these places are related [5], [6], [118], [119]. Topological information can be used to both increase the number of correct place matches and filter out incorrect matches [14], [84]. A probabilistic system like FAB-MAP can be run as a pure image retrieval process by assuming a uniform location prior at all steps, but performance improves when transition information is included through Bayesian filtering or similar techniques.

While image retrieval techniques can use an inverted index to improve efficiency, topological maps can use a location prior to speed up matching: the place recognition system only has to search places known to be close to the robot’s current position. A sampling-based method such as a particle filter can be used to sample possible places [12], [13], [111], [145]. The particles are resampled according to which places are the most likely, and can stay close by the current robot location if it is well-localized, or spread out across the whole environment if the robot is lost. Computation time is thus proportional to the number of particles, not the size of the environment [146].

Alternatively, since the number of loop closures in an environment is naturally sparse, Latif, Huang et al. [19] use topological information to formulate place recognition as a sparse convex L1-minimization problem, and apply efficient homotopy methods [147] to provide loop closure hypotheses.

The addition of topological information into the recognition process allows place recognition using low-resolution data and thus lower memory requirements. Using the sparse convex L1-minimization formulation, successful place recognition was achieved using images as small as 48 pixels [19]. Even in challenging scenarios where images are blurred or observed under different environment conditions such as different times of day, the use of topological information allows visual place

recognition using as few as 32 4-bit pixels per image, [148], [149].

C. Topological-metric maps

As image retrieval can be enhanced by adding topological information, topological maps can be enhanced by including metric information – distance, direction, or both – on the map edges. For example, both FAB-MAP [6] and SeqSLAM [118] are originally purely topological systems, but the addition of odometry information has been demonstrated to improve each system’s place recognition performance by CAT-SLAM [13] and SMART [150] respectively.

These topological-metric maps can be appearance-based, in which case metric information is only included as relative poses *between* each place node [151]–[154]. However, metric information about the position of landmarks or objects in a place can also be stored *within* each node [1], [2], [26], [141], [155]–[158]. The metric information within the topological place node can be stored as a sparse landmark map [2], [7], [76], or as a dense occupancy grid map [135] if depth information is extracted from the image data. Although the notion of dense spatial modeling using a truncated signed distance function (TSDF) representation can be traced back to the work of Moravec and Elfes [39] in the mid-1980s, it has become feasible only in the past few years, with the advent of GPU technology [135].

VI. RECOGNIZING PLACES: THE BELIEF GENERATION MODULE

Ultimately the purpose of a place recognition system is to determine whether a place has been seen before. Thus the central goal of any place recognition system is reconciling visual input with the stored map data to generate a belief distribution. This distribution provides a measure of likelihood or confidence that the current visual input matches a particular location in the robot’s map representation of the world. There is a general understanding that if two places descriptions appear similar there is a greater likelihood of them being captured at the same physical location, but the degree to which this is true depends on the particular environment. For example, repetitive environments may exhibit perceptual aliasing where different places look indistinguishable. Conversely, changing conditions may cause the same place to appear drastically different at different times.

A. Place recognition and SLAM

Place recognition plays an important role in pose graph SLAM algorithms by providing loop closure candidates [159]. Pose graphs, also known as view-based representations [160], [161], are widely utilized in modern SLAM systems because of their computational efficiency for fixed size maps, although they can suffer from an increase in computational requirements for long duration missions. Loop closure is vital for consistent mapping as it allows the system to correct drift in local odometry measurements [162], [163]. It can be decoupled from the online local update step, and many systems independently perform both SLAM-like local metric correction and topological-like loop closure [1], [2], [80],

[163]: a system can perform local metric correction using laser scan data [80], [163] or visual odometry [1], [2] while a separate global process looks for matches in order to close large loops.

If the place descriptions are appearance-based, and do not contain any metric information, but the map contains metric distances between places, the system can still use the loop closures to perform metric correction at the place level [151]–[154]. However, if the place descriptions contain metric information associated with the image features, as is the case for FrameSLAM [2], then a more precise correction can be performed. Maps that are purely topological or pure image retrieval do not provide any metric pose correction. In these cases, localization at a topological level occurs; that is, the system simply identifies the most likely location.

The place recognition maps that contain metric information both within and between the place descriptions can be used to perform a full metric SLAM solution. There are a wide range of SLAM techniques available as summarized in [164]–[166]. Thrun and Leonard [166] identify three key SLAM paradigms: Extended Kalman Filters (EKF) [37], [38], [167]–[169] and Rao-Blackwellized particle filters [170], as well as the pose graph approach discussed above [162], [163], [171]–[173]. Vision-based systems utilize all these methods: MonoSLAM [7] uses an EKF, while Rao-Blackwellized particle filters are used in [12], [174], [175] and pose graph optimization techniques in [2], [176].

B. Topological place recognition

If multiple streams of data are available a voting scheme [3], [5], [79], [96], [177] can be used. Ulrich and Nourbakhsh [5] used a Jeffrey divergence to compare color histograms and each color band votes for what it considers the most likely location. Depending on the votes, the system can be confident if the confident bands are unanimous and the total confidence is above a certain threshold, uncertain if none of the bands are sufficiently confident, or the total confidence value is too low, or confused if the confident bands disagree on the location.

If a system uses the bag-of-words model, inspired as it is by text-based document analysis, it may use the related Term Frequency-Inverse Document Frequency (TF-IDF) score [56], [114], [178]. Each visual word in an image has a TF-IDF score, which is made up of two parts – the term frequency, which measures how often the word appears in the image, and the inverse document frequency, which measures whether the word is common across all images. The TF-IDF score is then the product of these two values.

A probabilistic calculation can also be used to compute place matching likelihood, using a calculation based on Bayes theorem. Early examples of appearance-based probabilistic localization used Gaussians to represent probability [179], or a mixture of Gaussians combined with Expectation Maximization (EM) [180], or a Gaussian kernel [181] with Parzen smoothing [125]. Other choices for the observation likelihood include the use of TF-IDF for the observation likelihood, if a bag-of-words model is being used [83], [182]. Siagian and Itti [111], [183] use Monte Carlo Localization

(MCL) with two observation update steps each with an independent observation likelihood, one based on the segment likelihood and one based on the object likelihood. Garcia-Fidalgo and Ortiz [184] use the observation likelihood that relates the number of feature matches between two images to the overall number of features in the image, scaled by a normalizing constant.

The observational likelihood can also be computed via a data-driven approach. FAB-MAP [6], [87] is a probabilistic appearance-based localization system that uses a data-driven approach to calculating an observational likelihood. FAB-MAP uses a bag-of-words model with SIFT or SURF features for image description and calculates the distinctiveness of each word during a training phase. As a bag-of-words model may have many words – FAB-MAP has been used with a 100,000-word vocabulary [87] – the full joint probability distribution of the observed words (Fig. 8(a)) can be approximated by a naïve Bayes assumption (Fig. 8(b)) or a Chow-Liu tree [185] (Fig. 8(c)).

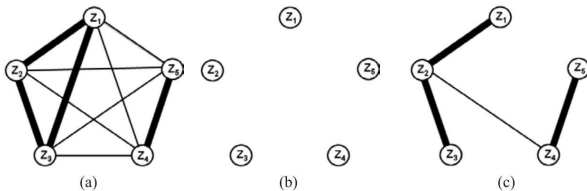


Fig. 8. FAB-MAP learns a probabilistic model of the relationship between word appearance and place recognition. (a) A full joint distribution takes into account the relationships between words (the thick lines between words represent those with the largest mutual information). (b) A naïve Bayes approximation of the full joint distribution ignores the mutual information between the words and assumes that all words appear independently. (c) A Chow-Liu tree approximates the full joint distribution as a junction tree where each word depends only on one other word (from [6]).

FAB-MAP handles the perceptual aliasing problem by considering not only whether two locations are similar in the sense that they have many visual words in common, but also whether the words in common are sufficiently rare that the locations can be considered distinctive. As a result, if two locations look similar but the words that appear are frequently observed, FAB-MAP will generate a low matching probability. FAB-MAP achieves this by using the denominator as a normalizing constant that is calculated over the set of all previously seen locations *and* the set of all locations that have not yet been visited.

Originally, the set of unvisited locations was modelled by randomly sampling from the Chow-Liu tree, and the probability that the robot was at a location that has not yet been observed was a user-defined parameter. However, Paul and Newman [60], [62], [186] presented an iterative learning mechanism to generate a representative set of the true distribution of the appearance of the world. Latent Dirichlet Allocation (LDA) [187] was used to cluster images into major topics that summarize how the world, as seen so far by the robot, appears. These topics are used to generate a sampling set that is proportional to what is common in the world – for example, foliage occurs frequently in many environments so should not be considered distinctive. The system learns

incrementally: after each deployment a better sampling set is created as the system incrementally learns about the world. Furthermore, an online-offline learning process is proposed – during the robot’s “down-time” further relevant data can be searched for on the internet to learn more about the world.

Olson [188] observes that “correct hypotheses generally agree with each other, whereas incorrect hypotheses tend to disagree with each other”. This property can be used to eliminate false positive matches by calculating a pair-wise consistency matrix between possible hypotheses and finding the most consistent set of hypotheses from the dominant eigenvectors. The same paper also observes that the amount of information required to generate a belief match should scale with the robot’s positional uncertainty. The system ensures this by requiring that local hypothesis matches cover a large physical space in comparison to the robot’s positional uncertainty, to ensure that the robot will not be incorrectly located within its uncertainty ellipse.

This approach contrasts with FAB-MAP’s requirement of a few highly distinctive matches. Instead, many matches are required over a large area, but these matches do not need to be particularly distinctive, as the geometrical relationship between the matches ensures the uniqueness of the hypothesis.

Biologically-inspired methods for place recognition mimic the known place cells structure in the rat hippocampus [116], [189]. In RatSLAM [116], a type of neural network known as a continuous attractor network (CAN) is used to model place cells (see Fig. 9). A continuous attractor network uses a combination of local excitation and global inhibition combined with input from ego-motion and visual sensors to perform localization. In a similar manner Giovannangeli, Gaussier et al. [189] use a place cell model to perform vision-based navigation in indoor and outdoor environments without a metric map.

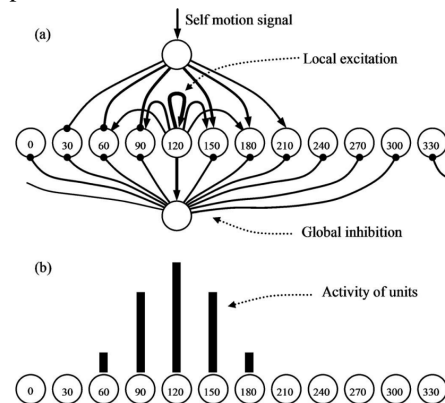


Fig. 9. Continuous attractor networks (CANs) are a type of neural network that can be used to model the behavior of place cells, head direction cells, and grid cells. (a) shows an example of a CAN used to model head direction cells. Each cell excites itself and units near itself (see local excitation arrows) and inhibits other cells. (b) shows a stable activity packet centered at 120° generated by the combination of local excitation and global inhibition with input from a motion input (from [116]).

C. Evaluation of place recognition systems

Topological place recognition systems are typically evaluated using precision and recall metrics and their relationship via a precision-recall curve. A system selects

matches based on a particular confidence measure. The correct matches are known as true positives, the incorrect matches are false positives, and matches that the system erroneously discards are false negative matches. Precision is defined as the proportion of selected matches that are true positive matches, and recall is the proportion of true positives to the total number of correct values, that is:

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

A perfect system would be one that achieves precision of 100% and recall of 100%. Precision and recall are often related to each other via a precision-recall curve which plots recall against precision for a range of confidence values.

Until recently, place recognition prioritized avoidance of false positive matches [6], as introducing false matches into a map could cause catastrophic failure. As a result, recall at 100% precision was the key metric for place recognition success. However, several methods for using topological information to correct false positive matches have been proposed [190]–[192] and attention has turned from eliminating all false positives to finding many potential place matches and then correcting any mismatches in a topological post-processing step. Increasing the number of potential matches is particularly important when performing place recognition in changing environments, when strict matching methods are liable to fail.

Furthermore, as place recognition systems transition from “demonstration” (typically with pre-recorded data sets) to “deployment” (operating in real-time on autonomous vehicles), the performance evaluation methodology may change further to include a consideration of the spatial distribution of place matches within the environment. For example, McManus, Churchill et al. [193] used the probability of travelling a given distance without a successful match as a measure of place recognition success. This metric expresses how evenly distributed the place matches are across the environment and is an important measure for the overall integrity of a navigation system that uses place recognition as a module.

VII. VISUAL PLACE RECOGNITION IN CHANGING ENVIRONMENTS

Early place recognition systems often implicitly used the simplifying assumption that the visual appearance of each place would not change over the course of the experiment. However, as robotic systems operate in ever-larger, uncontrolled environments and for longer time periods, it has rapidly become apparent that this assumption is no longer valid. Consequently, in recent years there has been a growing focus on creating persistent robotic navigation systems, including persistent place recognition techniques. The ability to localize in and generate maps of dynamic environments has been identified as being of key importance [194]. This section revisits each of the previous concepts – how a place can be

represented, how the mapping frameworks work and how the belief generation process works – and discusses how each has to change to manage a changing environment.

A. Describing places in changing environments

It is clear that the appearance of a place can vary greatly over time due to a large number of causes including changes in lighting and weather (see Fig. 1). There are two methods for performing place recognition when faced with appearance change – the first tries to find a condition-invariant description of the place, the way local feature descriptors are designed to be scale-, rotation- and illumination-invariant. The second method tries to learn how appearance change occurs.

1) Invariant methods

The difficulty of matching places in changing environments using conventional local features is a significant one for persistent robot navigation: Furgale and Barfoot [117] observed that the non-repeatability of SURF features due to changing appearance, particularly lighting change, was a major cause of failure during visual-teach-and-repeat experiments. Existing image description methods have been tested to determine their robustness to illumination and other change. In [195], Valgren and Lilienthal tested SIFT features and a number of SURF variants across change in lighting, cloud cover, and seasonal conditions. The SURF variants all outperformed SIFT, but none of the tested features were found to be robust across all conditions. However, in later work [100] the authors combined U-SURF [73], the most successful SURF variant, with a consistency check using the epipolar constraint, and achieved between 80% and 100% correct matching within small (40 image) datasets.

Ross, English et al. [196], [197] studied the effect of lighting change on features using time-lapse footage across full days to determine the illumination sensitivity of each descriptor. The feature keypoints were predefined within each image, and only the variance of the feature descriptor was tested, in contrast to the work of Valgren and Lilienthal [100], [195] which tested the combined effect of feature detector and descriptor. The U-SIFT [72] descriptor was shown to display the greatest lighting invariance of the tested descriptors.

Instead of using point features such as SIFT or SURF, other descriptors can be chosen. Whole-image descriptors have been used in systems such as SeqSLAM [118], [119], [198] that demonstrate robustness against environmental change. However, as for other description methods, too drastic a change in appearance will cause system failure [111] and whole-image descriptors also suffer from the additional problem of sensitivity to viewpoint change [199]. Edge features can be used in appropriate environments [174], [200], as they are invariant to lighting, orientation and scale [200]. Nuske, Roberts et al. [200] used line-based localization to localize against an existing map with a fish-eye camera and tested it in an outdoor industrial area under various lighting conditions across times of day from 7:00 to 17:00. Borges, Zlot et al. [201] extended this system to generate its own edge map using 3D laser data for localization. However data association using edge features can be challenging [174].

Techniques such as shadow removal [202] and the use of an illumination invariant color space [193] can lessen the effect of appearance variability caused by illumination change. Alternatively, a hardware-based solution to place recognition in variable lighting conditions can be used. McManus, Furgale et al. [203] used scanning laser-rangefinders to create “camera-like” images that were not affected by the illumination of the scene. This solution had the advantage of being applicable in complete darkness. A long-wave infrared thermal imaging camera is another sensor that can be deployed in a manner similar to a standard camera but which responds differently to lighting variance. Maddern and Vidas [204] showed thermal imaging cameras can provide improved place recognition at night-time when visible light cameras fail.

Convolutional Neural Networks (CNNs) have recently been used as robust feature extractors for place recognition in changing environments. Exploring the utility of CNNs for place recognition has been motivated by their ability to learn generic features that are transferrable to a variety of related but different visual tasks [205], [206]. [207], [208] utilized CNN features as holistic image descriptors and analyzed the robustness of different layers against visual appearance and viewpoint changes. They concluded that mid-level features exhibit a robustness against appearance changes, while higher level features are more robust against changes in viewpoint and carry more semantic information that can be used to partition the search space [208].

One aspect of visual data that has not been investigated in depth for changing environments is that of color. While conventional images descriptors such as SURF and BRISK operate on grayscale images, most available cameras capture color images, which have the potential to provide new and interesting information about place recognition in changing environments. Color information presents an interesting paradox for place recognition in changing environments: it is known to perform poorly as a feature when the illumination of a scene changes [196], but conversely, relative color information contains information about lighting that can improve place recognition dramatically by identifying and removing shadows [202]. Illumination invariant images use relative color information and are more reliable for place recognition during the day, but are out-performed by color images at night, when the underlying assumptions about black-body illumination are violated [209].

2) *Learning methods*

The alternative to invariant approaches is to learn a relationship between how places appear at different times. These methods assume that places change appearance in a similar way across an environment, and so change learned during training can be generalized to previously unseen locations. This assumption has been tested by observing static webcams from different locations [210], [211] and demonstrating that the most significant transformations across time are similar across different places. Furthermore, a training set of locations can be used to compute a principal component basis that encodes new locations with only a small loss of accuracy.

Ranganathan, Matsumoto et al. [212] learned a fine vocabulary [213]; a fine vocabulary is similar to a bag-of-words model in that it segments a descriptor space, such as SIFT descriptors, but it does so very finely – into over 16 million words in [213]. The system then learned a probability distribution over these words. The motivation for the fine vocabulary is the observation that descriptors transform in a highly non-linear way due to illumination change, changing viewpoint and other effects, and learning a distribution of alternative words allows these changes to be learned and quantified. In [212] the distribution was learned over multiple training runs over the same environment and features were matched across different illumination conditions to generate the probability distribution. Improved performance was reported over using a conventional vocabulary tree [95], with an additional 10%-15% of the dataset being correctly matched. The distance metric was also compared and the symmetric KL-divergence was shown to out-perform either the standard descriptor distance metric or a probability distance metric.

Using webcam footage, Carlevaris-Bianco and Eustice [214] tracked image patches over different lighting conditions to generate a large set (3 million features) of positive and negative examples. From this data, a neural network learning technique [215] mapped the patches into a new space in which positive matches were close together, according to the Euclidean distance, and negative matches were further away. The mapped descriptors were shown to be substantially more successful at place recognition than SIFT and SURF descriptors – compared to SURF descriptors, an additional 10% of the test locations were correctly matched.

Neubert, Sünderhauf et al. [18] learned a visual translation between two different seasons. Training images from two different seasons were segmented using SLIC superpixels [216]. The superpixels were described using a color histogram and a SURF descriptor, and a dictionary of translations of superpixels from one season to another season was learned. Similarly, Lowry, Milford et al. [217] learned a linear transformation from images captured in the morning to images captured in the late afternoon. However, for such appearance translation to be successful, the pairs of training images must be well aligned.

Learning-based methods frequently require a supervised training phase, which implies that the likely appearance change is known and that relevant training data is available. Lowry, Wyeth et al. [218] proposed an unsupervised learning method for place recognition in changing environments. Instead of attempting to predict the appearance of a location, the system instead identified and removed potentially changing aspects of each observation.

B. *Remembering places in changing environments*

If the environment is changing, the map also needs to change to continue representing the environment. The system must determine what to remember and what to forget. It may also be beneficial for the system to maintain multiple representations of a place, as places can vary between different configurations. This section presents mapping

frameworks for place recognition that have the capacity to handle changing environments in one of these two ways – either by deciding what to remember and what to forget, and/or by remembering multiple different representations. These systems are not all specific to vision-based systems, and many have been designed to handle laser data, but demonstrate concepts that are relevant to any sensor modality or map framework.

1) *Remembering and forgetting data*

In a dynamic environment, each place representation must be updated as new observations are obtained by the robot. A balance has to be found between using recent observations to overwrite obsolete information, and not allowing fleeting events to overwrite the status quo. However, it is difficult to determine which events are transient and which are worth remembering. Drawing inspiration from concepts in neuroscience, Biber and Duckett [219] referred to this as the “stability-plasticity dilemma”. Biological brains can inspire solutions for coping with this dilemma: concepts such as sensory memory, short-term memory and long-term memory found in human memory models have been co-opted to create decision models for remembering and forgetting.

One biologically inspired mapping system passes sensor information through an analogue of sensory memory to short-term memory and long-term memory storage areas [220], [221]. In the first stage, a selective attention mechanism decides which information will be upgraded from sensory memory to short-term memory, based on information from the long-term memory. The second stage involves using a rehearsal mechanism to determine which information will be transferred from short-term to long-term memory. Using attention and rehearsal mechanisms ensures that more persistent, stable and frequently occurring features are remembered, whilst transient features are forgotten. Elements must be seen and recognized sufficiently often before they are considered for promotion to a higher level of memory. Furthermore, obsolete features are slowly filtered out of the long-term memory. There is a complementary problem of which elements to ‘remember’, which typically uses similar criteria [220], [222] to the forgetting process.

Andrade-Cetto and Sanfeliu [223] required that features be trustworthy and reliable as well as up-to-date in order to be retained, while Bailey [222] considered a usefulness criteria based on visibility – a feature that can be blocked by other elements of the environment is liable to suffer from occlusion errors and be less useful in the future. Johns and Yang [102] and Hafez, Singh et al. [224] used a bag-of-words model and applied a quality measure to determine useful features to retain, considering both feature distinctiveness and feature reliability when generating a model of a location. Johns and Yang [225] also proposed a generative bag-of-words model that considered the variance as well as the mean value of each data point when matching scenes.

2) *Multiple representations of the environment*

Not only do places change in appearance over time, but they may also change in a cyclic manner that cannot be represented by a single description. During a two-week office-based

experiment [226], Milford and Wyeth noted that “the weakness is that the system deals rather inefficiently with cyclic changes such as day–night time cycles. Over a full night of operation, the pruning process gradually develops the experience map representation into one suited to localization at night time, somewhat hindering localization in the morning.” These observations were corroborated by Ranganathan, Matsumoto et al. [212], who stated that for an indoor office environment, consistently good localization through the 24-hour cycle would require around 3-4 images per location. Rather than continuously remembering and forgetting information, the map should hold multiple representations of the area – whether at a place or higher level.

A place recognition system can use multiple maps of the same environment. In the work of Biber and Duckett, each map remembered a different timescale [227]. Some of these maps represented short-term memory and were updated frequently whilst others were analogous to long-term memory and are not updated for hours, days, or weeks. Keeping maps that updated at different timescales ensured that old mapping data was not immediately overwritten by a temporary change in the environment. Instead the most static elements were reinforced over time, whilst transient events were filtered out. Place recognition was performed by selecting the local map that best fitted the current sensor data.

Systems that maintain multiple maps of the same environment may also add new map configurations only when they are necessary, rather than according to a pre-set timeframe [221]. Furthermore, Stachniss and Burgard [228] noted that not every place needs multiple representations – certain areas such as doorways may exhibit more change than the rest of the environment. Such areas may only possess a few key configurations – for example, a door may be open or closed – so the world can be described sufficiently accurately using a finite number of submaps. Each region in which dynamic activity is observed was segmented from the rest of the map in a submap. Fuzzy k -means clustering was used with the Bayesian Information Criterion to determine the optimal number of typical configurations of this area. Using submaps to segregate dynamic areas allowed multiple environmental configurations where necessary whilst keeping the map manageable.



Fig. 10. The varying appearance of a changing environment may require a system store multiple representations of each place. This image (from [229]) shows the number of robot “experiences” stored during repeated traversals of a path over a number of months. While most places require 5-10 experiences (shown in blue) some regions require as many as 30 (shown in red).

Elements of a scene that are moving when the robot observes them must be detected and may also be removed [230], [231]. However, there are often semi-static elements that are not obviously moving but appear and disappear over time. While these elements can simply be removed as unreliable [69], [232] it is also possible that such elements may be temporarily useful for localization in specific parts of an environment [233]. For example, in a car park building the static elements such as the walls can be far away and not particularly distinctive, and so are not useful for localization while the semi-static parked cars are many and relatively distinctive, and can be used for localization for a matter of hours or a day, before being forgotten and replaced. If this is the case, temporary maps are created when the robot observations do not match the expected results of the provided static map. The temporary maps are discarded when they fail to adequately match the robot observations over multiple consecutive time steps.

The systems presented above [221], [227], [228], [233] were designed for metric systems. Multiple representations can also be generated for appearance-based systems if multiple training runs are available. Johns and Yang [102] used feature co-occurrence maps generated during five training runs on a 20 km urban road-based dataset between 14:00 and 22:00. Localization can then be achieved on the same route at times interpolated between the five runs.

McManus, Upcroft et al. [120] used multiple training runs through an environment to learn scene signatures – locally distinctive elements of a place that are also stable over changes in appearance. For each location within the environment, image patches are selected that specifically demonstrate both distinctiveness and stability. The selected patches were described using HOG descriptors [121], and used to train an SVM classifier for each location. Using scene signatures for each places allowed 100% correct place recognition in a 31 location dataset, while SURF features performed poorly, particular in rainy and foggy conditions.

If the appearance of the environment is assumed to be affected by a series of hidden periodic processes, spectral analysis such as Fourier analysis can be used to predict the most likely appearance of a location from multiple training passes at a particular time in the future. Krajník, Fentanes et al. [234] learned and modeled these processes over an environment and demonstrated that this information can halve the number of place recognition errors when localizing three months later.

All of the systems described above share an underlying assumption – that the robot knows where it is sufficiently well to match different representations of the same location together, even if the representations are visually dissimilar. A map cannot be updated if the system does not know which location to update and, in a changing environment, it may not

be possible to know exactly where the robot is. To avoid this assumption, Churchill and Newman proposed a plastic map formulation [15] that explicitly localizes within robot “experiences” rather than physical locations. A new experience is generated each time a robot visits a location that it does not recognize, and the map may implicitly have multiple representations of each location, depending on the difficulty of matching at that particular location (see Fig. 10). However, unlike the systems discussed previously, the multiple representations will not necessarily be linked together as the same physical place. The plastic map is more informative if the system can recognize and link more experiences together. However, it is a pragmatic approach that allows for graceful place recognition failure without catastrophic map collapse.

Retaining multiple representations of each location increases the place recognition search space and can decrease efficiency unless only a subset of representations is used for comparison. Because observations captured at similar times tend to demonstrate similar appearance characteristics, future potential matches can be probabilistically selected based on the system’s current localization belief. Carlevaris-Bianco and Eustice [235] approximated the likelihood of two location exemplars being “co-observed” within a short time-frame with a Chow-Liu tree, while Linegar, Churchill et al. [236] used “path memory” to select past experiences as candidate matches and improve place recognition without increasing computation time.

C. Recognizing places in changing environments

Integrating appearance change into a place recognition system requires some key alterations to the belief generation process. Firstly, as discussed above, changing environments require multiple representations of each place. If this is the case, a system may select the best map given its current sensor data [227] or it may try to predict the most likely appearance matches [18], [234]–[236].

Alternatively, the place recognition system may run multiple hypotheses in parallel. Churchill and Newman [15] assigned every saved experience its own localizer that reports whether or not the robot is successfully localized within that environment, while Morris, Dayoub et al. [221] performed filtering over possible map configurations as well as possible robot poses. Instead of selecting the single map that best matches the current sensor data, the system instead actively tracks the N best navigation hypotheses in multiple maps, while pending hypotheses are maintained and swapped out when an active hypothesis drops below the best pending hypothesis. Using multiple map hypotheses was reported to decrease the mean path error in an indoor office experiment by as much as 80%.

One factor for place recognition in changing environments is that topological information becomes more important as incoming sensor data becomes less reliable and more difficult to match to previous observations [118], [119]. It has been observed that matching image sequences rather than individual images can improve place recognition in general, and

particularly in changing environments [14], [84], [118], [149], and image sequences can be integrated with conditional random fields [237] to identify and if necessary verify loop closures [14].

The place recognition systems that are most successful in changing environments exploit the assumption that the system is not just passing through a particular place, but traversing the same or a very similar path through the environment. SeqSLAM [118] demonstrated that image sequences can perform place recognition in particularly visually challenging environments. The original version assumed a similar velocity profile between traversals. Methods to deal with this limitation include searching non-linear paths as well as linear paths [102] through the image similarity matrix and using odometry input to linearize the signal [150]. Liu and Zhang [238] used a particle filter to improve the computation efficiency over the exhaustive search process and achieved a 10 times speed-up factor with equivalent performance at 100% precision.

Naseer, Spinello et al. [119] exploited sequence information by formulating image matching as a minimum cost flow. Flow networks are directed graphs with a source node and a sink node, which for path-based place recognition represent the start of the traversal and the end of the traversal respectively. By equating image comparison values to flow cost, the formulation found the optimal sequence through the environment. Differing velocity profiles were handled by allowing nodes to be either matching or hidden. Similarly, Hansen and Browning [239] used Hidden Markov Models to determine the most likely path through an environment using the Viterbi algorithm.

VIII. CONCLUSION

Visual place recognition has made great advances in the last 15 years, but we are still a long way from a universal place recognition system for robots that is robust and widely applicable across a range of robotic platforms and varying environments. Here we highlight several promising avenues of ongoing and future research that are moving us closer towards this outcome.

The most successful approaches to combatting changing appearance typically do so at the cost of viewpoint invariance or increased training requirements. As discussed above, as sensor information becomes less reliable, it can be compensated for by topological information, which requires not only viewpoint invariance at a single point, but along a possibly quite long path. Some potential avenues include using image patches rather than whole images, as image patches have much of the condition invariant advantages of whole images while allowing some coarse viewpoint invariance, and investigating the use of deep learning features which also have some viewpoint invariant characteristics.

Visual place recognition is benefitting from research in other fields, particularly the great strides being achieved in computer vision in the fields of deep learning, image classification, object recognition, video description. While techniques such as convolutional neural networks depend on Big Data and Big Compute, techniques such as cloud robotics

and online / offline processing paradigms could be exploited even by small, cheap mobile platforms. Developments in GPU hardware and novel camera sensors will inspire new concepts in place recognition as well as improving the efficiency and robustness of existing approaches.

Research in place recognition can also benefit from the ongoing research in object detection and scene classification. By exploiting object detections, it is possible to learn that objects such as buildings are useful for long-term place recognition, objects such as pedestrians should be ignored, and objects such as cars might be useful depending on the semantic and temporal context. An increased robustness to structural changes can be achieved by exploiting knowledge about which objects are dynamic or static and how that property is depending on the temporal and semantic context – for example, cars in a parking garage can temporarily provide useful place recognition cues. Exploiting the expressiveness of convolutional neural networks by training or fine-tuning such networks specifically for the task of place recognition is a worthwhile direction for future research.

Visual place recognition systems can also exploit context. Although places change drastically in appearance, the relative location of places remains unchanged. This fact is integrated into belief generation modules by using location priors, recursive filtering and path-based sequences of images, and the dependence on these techniques increases as the variation in the visual appearance of the environment increases. The use of other sources of contextual information also has the potential to improve place recognition capability – knowledge about the time of day, or the current weather conditions can also change how the place recognition system interprets the incoming visual data.

Semantic scene context can furthermore limit the search space for place recognition to semantically similar scenes to ensure scalability towards long-term autonomy. Semantic context can support learning and predicting the changes in a scene and help increase the robustness against environmental condition changes. Semantic mapping also has the potential to reduce memory requirements – imagine a house map only requiring words such as “kitchen”, “bedroom”, and “bathroom” to describe places – and current research in topic modeling, coresets and other semantic compression methods is already showing promise, as is the use of objects as high-level place recognition features.

Finally, what can visual place recognition offer to other research tasks? By necessity and opportunity, visual place recognition has taken up the challenge to solve condition invariant recognition to a degree that many fields have not, albeit under a more tightly constrained task specification than other tasks such as scene interpretation. The experience gained in developing robust features, in addressing the combination of both appearance change and viewpoint change and other challenges may have valuable applications both in other robotic tasks such as object recognition and classification in the wild, and a diverse range of other areas including remote sensing, environmental monitoring and tasks that require recognition and identification in uncontrolled environments.

REFERENCES

- [1] G. Sibley, C. Mei, I. Reid, and P. Newman, "Vast-scale outdoor navigation using adaptive relative bundle adjustment," *Int. J. Rob. Res.*, vol. 29, no. 8, pp. 958–980, 2010.
- [2] K. Konolige and M. Agrawal, "FrameSLAM: From bundle adjustment to real-time visual mapping," *IEEE Trans. Robot.*, vol. 24, no. 5, pp. 1066–1077, 2008.
- [3] G. Schindler, M. Brown, and R. Szeliski, "City-scale location recognition," in *Computer Vision and Pattern Recognition (CVPR 2007)*, *IEEE Conference on*, 2007, pp. 1–7.
- [4] M. Milford, G. Wyeth, and D. Prasser, "RatSLAM: A hippocampal model for simultaneous localization and mapping," *Robot. Autom. (ICRA 2004)*, *2004 IEEE Int. Conf.*, pp. 403–408, 2004.
- [5] I. Ulrich and I. Nourbakhsh, "Appearance-based place recognition for topological localization," in *Robotics and Automation (ICRA 2000)*, *2000 IEEE International Conference on*, 2000, vol. 2, pp. 1023–1029.
- [6] M. Cummins and P. Newman, "FAB-MAP: Probabilistic localization and mapping in the space of appearance," *Int. J. Rob. Res.*, vol. 27, no. 6, pp. 647–665, 2008.
- [7] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, "MonoSLAM: Real-time single camera SLAM," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 6, pp. 1052–1067, 2007.
- [8] C. G. Harris and J. M. Pike, "3D positional integration from image sequences," *Image Vis. Comput.*, vol. 6, no. 2, pp. 87–90, 1988.
- [9] J. Neira, M. I. Ribeiro, and J. D. Tardós, "Mobile robot localization and map building using monocular vision," in *In The 5th Symposium for Intelligent Robotics Systems*, 1997.
- [10] M. Bosse, R. Rikoski, J. Leonard, and S. Teller, "Vanishing points and three-dimensional lines from omni-directional video," in *IEEE International Conference on Image Processing (ICIP)*, 2002, vol. 19, no. 6, pp. 417–430.
- [11] R. Eustice, H. Singh, J. J. Leonard, M. Walter, and R. Ballard, "Visually Navigating the RMS Titanic with SLAM Information Filters," in *Robotics: Science and Systems*, 2005, pp. 57–64.
- [12] R. Sim, P. Elinas, M. Griffin, and J. J. Little, "Vision-based SLAM using the Rao-Blackwellised particle filter," in *IJCAI Workshop on Reasoning with Uncertainty in Robotics*, 2005, vol. 14, no. 1, pp. 9–16.
- [13] W. Maddern, M. Milford, and G. Wyeth, "CAT-SLAM: Probabilistic localisation and mapping using a continuous appearance-based trajectory," *Int. J. Rob. Res.*, vol. 31, no. 4, pp. 429–451, 2012.
- [14] C. Cadena, D. Gálvez-López, J. D. Tardós, and J. Neira, "Robust place recognition with stereo sequences," *Robot. IEEE Trans.*, vol. 28, no. 4, pp. 871–885, 2012.
- [15] W. Churchill and P. Newman, "Experience-based navigation for long-term localisation," *Int. J. Rob. Res.*, vol. 32, no. 14, pp. 1645–1661, 2013.
- [16] F. Dayoub, T. Morris, B. Upcroft, and P. Corke, "Vision-only autonomous navigation using topometric maps," in *IEEE/RSJ International Conference on Intelligent Robots and Systems November 3-8, 2013 at Tokyo Big Sight, Japan*, 2013.
- [17] A. Murillo, G. Singh, J. Košecák, and J. Guerrero, "Localization in urban environments using a panoramic Gist descriptor," *IEEE Trans. Robot.*, vol. 29, no. 1, pp. 146–160, 2013.
- [18] P. Neubert, N. Sünderhauf, and P. Protzel, "Superpixel-based appearance change prediction for long-term navigation across seasons," *Rob. Autom. Syst.*, 2014.
- [19] Y. Latif, G. Huang, J. Leonard, and J. Neira, "An online sparsity-cognizant loop-closure algorithm for visual navigation," in *Proceedings of Robotics: Science and Systems Conference (RSS)*, 2014.
- [20] E. C. Tolman, "Cognitive maps in rats and men," *Psychol. Rev.*, vol. 55, no. 4, p. 189, 1948.
- [21] A. T. Bennett, "Do animals have cognitive maps?," *J Exp Biol*, vol. 199, no. Pt 1, pp. 219–224, 1996.
- [22] R. Jensen, "Behaviorism, latent learning, and cognitive maps: needed revisions in introductory psychology textbooks," *Behav Anal.*, vol. 29, no. 2, pp. 187–209, 2006.
- [23] K. Lynch, *The Image of the City*. Cambridge, MA: MIT press, 1960.
- [24] B. Kuipers, "Modeling Spatial Knowledge," *Cogn. Sci.*, vol. 2, no. 2, pp. 129–153, 1978.
- [25] B. Kuipers and Y.-T. Byun, "A robot exploration and mapping strategy based on a semantic hierarchy of spatial representations," *Rob. Autom. Syst.*, vol. 8, no. 1, pp. 47–63, 1991.
- [26] B. Kuipers, "The spatial semantic hierarchy," *Artif. Intell.*, vol. 119, no. 1, pp. 191–233, 2000.
- [27] F. Strumwasser, "Long-term recording from single neurons in brain of unrestrained mammals," *Science (80-)*, vol. 127, no. 3296, pp. 469–470, 1958.
- [28] J. O'Keefe and J. Dostrovsky, "The hippocampus as a spatial map. Preliminary evidence from unit activity in the freely-moving rat," *Brain Res.*, vol. 34, no. 1, pp. 171–175, 1971.
- [29] J. O'Keefe, "Place units in the hippocampus of the freely moving rat," *Exp. Neurol.*, vol. 51, no. 1, pp. 78–109, 1976.
- [30] M. A. Wilson and B. L. McNaughton, "Dynamics of the hippocampal ensemble code for space," *Science (80-)*, vol. 261, no. 5124, pp. 1055–1058, 1993.
- [31] J. O'Keefe and D. H. Conway, "Hippocampal place units in the freely moving rat: why they fire where they fire," *Exp. Brain Res.*, vol. 31, no. 4, pp. 573–590, 1978.
- [32] J. S. Taube, R. U. Muller, and J. B. Ranck Jr., "Head-direction cells recorded from the postsubiculum in freely moving rats. II. Effects of environmental manipulations," *J Neurosci*, vol. 10, no. 2, pp. 436–447, 1990.
- [33] T. Hafting, M. Fyhn, S. Molden, M.-B. Moser, and E. I. Moser, "Microstructure of a spatial map in the entorhinal cortex," *Nature*, vol. 436, no. 7052, pp. 801–806, 2005.
- [34] E. I. Moser, E. Kropff, and M. B. Moser, "Place cells, grid cells, and the brain's spatial representation system," *Annu Rev Neurosci*, vol. 31, pp. 69–89, 2008.
- [35] K. M. Gothard, W. E. Skaggs, and B. L. McNaughton, "Dynamics of mismatch correction in the hippocampal ensemble code for space: interaction between path integration and environmental cues," *J Neurosci*, vol. 16, no. 24, pp. 8027–8040, 1996.
- [36] A. D. Redish, E. S. Rosenzweig, J. D. Bohanick, B. L. McNaughton, and C. A. Barnes, "Dynamics of hippocampal ensemble activity realignment: time versus space," *J Neurosci*, vol. 20, no. 24, pp. 9298–9309, 2000.
- [37] R. Chatila and J. P. Laumond, "Position referencing and consistent world modeling for mobile robots," in *Robotics and Automation (ICRA 1985)*, *1985 IEEE International Conference on*, 1985, pp. 138–145.
- [38] R. Smith, M. Self, and P. Cheeseman, "A stochastic map for uncertain spatial relationships," in *Proceedings of the 4th International Symposium on Robotics Research*, 1988, pp. 467–474.
- [39] H. Moravec and A. E. Elfes, "High resolution maps from wide angle sonar," in *Robotics and Automation (ICRA 1985)*, *1985 IEEE International Conference on*, 1985, pp. 116–121.
- [40] R. A. Brooks, "Visual map making for a mobile robot," in *Robotics and Automation (ICRA 1985)*, *1985 IEEE International Conference on*, 1985, vol. 2, pp. 824–829.
- [41] H. Durrant-Whyte, D. Rye, and E. Nebot, "Localization of autonomous guided vehicles," in *Proceedings of the 8th International Symposium on Robotics Research*, 1995, pp. 613–625.
- [42] H. Shatkay and L. P. Kaelbling, "Learning geometrically-constrained hidden Markov models for robot navigation: bridging the geometrical-topological gap," *J. Artif. Intell. Res.*, 2002.
- [43] G. Dudek, M. Jenkin, E. Milios, and D. Wilkes, "Robotic exploration as graph construction," *IEEE Trans. Robot. Autom.*, vol. 7, no. 6, pp. 859–865, 1991.
- [44] E. Fabrizi and A. Saffiotti, "Extracting topology-based maps from gridmaps," *Robot. Autom. (ICRA 2000)*, *2000 IEEE Int. Conf.*, vol. 3, 2000.
- [45] T. Bailey, E. M. Nebot, J. K. Rosenblatt, and H. F. Durrant-Whyte, "Robust distinctive place recognition for topological maps," *Int. Conf. F. Serv. Robot.*, pp. 347–352, 1999.
- [46] P. Beeson, N. K. Jong, and B. Kuipers, "Towards autonomous topological place detection using the extended Voronoi graph," in *Robotics and Automation (ICRA 2005)*, *2005 IEEE International Conference on*, 2005, pp. 4373–4379.
- [47] D. Kortenkamp, L. D. Baker, and T. Weymouth, "Using gateways to build a route map," in *Intelligent Robots and Systems (IROS 1992)*, *1992 IEEE/RSJ International Conference on*, 1992, vol. 3, pp. 2209–2214.

- [48] A. Ranganathan, "PLISS: Detecting and Labeling Places Using Online Change-Point Detection," in *Robotics: Science and Systems*, 2010.
- [49] G. Tsechpenakis, D. N. Metaxas, C. Neidle, and O. Hadjiiladis, "Robust Online Change-point Detection in Video Sequences," in *Computer Vision and Pattern Recognition Workshop, 2006. CVPRW '06. Conference on*, 2006, p. 155.
- [50] L. Itti and P. Baldi, "A principled approach to detecting surprising events in video," in *Computer Vision and Pattern Recognition (CVPR 2005), IEEE Computer Society Conference on*, 2005, vol. 1, pp. 631–637 vol. 1.
- [51] S. R. Esterby and A. H. El-Shaarawi, "Inference about the point of change in a regression model," *Appl. Stat.*, pp. 277–285, 1981.
- [52] A. Ranganathan and F. Dellaert, "Bayesian surprise and landmark detection," in *Robotics and Automation, 2009. ICRA '09. IEEE International Conference on*, 2009, pp. 2017–2023.
- [53] H. Korrapati, J. Courbon, Y. Mezouar, and P. Martinet, "Image Sequence Partitioning for outdoor mapping," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, 2012, pp. 1650–1655.
- [54] A. Chapoulie, P. Rives, and D. Filliat, "Topological segmentation of indoors/outdoors sequences of spherical views," in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, 2012, pp. 4288–4295.
- [55] L. Murphy and G. Sibley, "Incremental unsupervised topological place discovery," in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, 2014, pp. 1312–1318.
- [56] T. Nicosevici and R. Garcia, "Automatic Visual Bag-of-Words for Online Robot Navigation and Mapping," *Robot. IEEE Trans.*, vol. 28, no. 4, pp. 886–898, 2012.
- [57] C. Tzu-Chuan and C. Meng Chang, "Using Incremental PLSI for Threshold-Resilient Online Event Analysis," *Knowl. Data Eng. IEEE Trans.*, vol. 20, no. 3, pp. 289–299, 2008.
- [58] M. Volkov, G. Rosman, D. Feldman, J. W. F. III, and D. Rus, "Coresets for Visual Summarization with Applications to Loop Closure," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2015.
- [59] D. Feldman, M. Schmidt, and C. Sohler, "Turning big data into tiny data: Constant-size coresets for k-means, PCA and projective clustering," in *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2013, New Orleans, Louisiana, USA, January 6-8, 2013*, 2013, pp. 1434–1453.
- [60] R. Paul, D. Rus, and P. Newman, "How was your day? Online visual workspace summaries using incremental clustering in topic space," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, 2012, pp. 4058–4065.
- [61] Y. Girdhar and G. Dudek, "Efficient on-line data summarization using extremum summaries," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2012.
- [62] R. Paul, D. Feldman, D. Rus, and P. Newman, "Visual precis generation using coresets," in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, 2014, pp. 1304–1311.
- [63] Y. Girdhar, P. Giguere, and G. Dudek, "Autonomous adaptive exploration using realtime online spatiotemporal topic modeling," *Int. J. Rob. Res.*, vol. 33, no. 4, pp. 645–657, 2014.
- [64] D. Filliat and J.-A. Meyer, "Map-based navigation in mobile robots: I. a review of localization strategies," *Cogn. Syst. Res.*, vol. 4, no. 4, pp. 243–282, 2003.
- [65] M. R. Walter, S. Hemachandra, B. Homberg, S. Tellex, and S. Teller, "Learning semantic maps from natural language descriptions," 2013.
- [66] A. Torralba, K. P. Murphy, W. T. Freeman, and M. A. Rubin, "Context-based vision system for place and object recognition," in *Computer Vision (ICCV 2003), Ninth IEEE International Conference on*, 2003, pp. 273–280.
- [67] D. Filliat, E. Battesti, S. Bazeille, G. Duceux, A. Gepperth, L. Harrath, I. Jebari, R. Pereira, A. Tapus, C. Meyer, I. Sio-Hoi, R. Benosman, E. Cizeron, J. C. Mamanna, and B. Pothier, "RGBD object recognition and visual texture classification for indoor semantic mapping," in *Technologies for Practical Robot Applications (TePRA), 2012 IEEE International Conference on*, 2012, pp. 127–132.
- [68] A. Pronobis and P. Jensfelt, "Large-scale semantic mapping and reasoning with heterogeneous modalities," in *Robotics and Automation (ICRA 2012), 2012 IEEE International Conference on*, 2012.
- [69] R. Biswas, B. Limketkai, S. Sanner, and S. Thrun, "Towards object mapping in non-stationary environments with mobile robots," in *Intelligent Robots and Systems (IROS 2002), 2002 IEEE/RSJ International Conference on*, 2002, pp. 1014–1019.
- [70] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. J. Kelly, and A. J. Davison, "SLAM++: Simultaneous localisation and mapping at the level of objects," in *2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013*, 2013, pp. 1352–1359.
- [71] J. Wu, H. I. Christensen, and J. M. Rehg, "Visual place categorization: Problem, dataset, and algorithm," in *Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on*, 2009, pp. 4763–4770.
- [72] D. Lowe, "Object recognition from local scale-invariant features," in *IEEE International Conference on Computer Vision (ICCV)*, 1999, vol. 2, pp. 1150–1157.
- [73] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, "Speeded-Up Robust Features (SURF)," *Comput. Vis. Image Underst.*, vol. 110, no. 3, pp. 346–359, 2008.
- [74] A. Oliva and A. Torralba, "Building the gist of a scene: The role of global image features in recognition," *Prog. Brain Res.*, vol. 155, pp. 23–36, 2006.
- [75] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *Int. J. Comput. Vis.*, vol. 42, no. 3, pp. 145–175, 2001.
- [76] S. Se, D. Lowe, and J. Little, "Mobile robot localization and mapping with uncertainty using scale-invariant visual landmarks," *Int. J. Rob. Res.*, vol. 21, no. 8, pp. 735–758, 2002.
- [77] H. Andreasson and T. Duckett, "Topological localization for mobile robots using omni-directional vision and local features," in *Proc. of the 5th IFAC Symposium on Intelligent Autonomous Vehicles (IAV)*, 2004.
- [78] E. Stumm, C. Mei, and S. Lacroix, "Probabilistic place recognition with covisibility maps," in *Intelligent Robots and Systems (IROS 2013), 2013 IEEE/RSJ International Conference on*, 2013, pp. 4158–4163.
- [79] J. Košecká, F. Li, and X. Yang, "Global localization and relative positioning based on scale-invariant keypoints," *Rob. Auton. Syst.*, vol. 52, no. 1, pp. 27–38, 2005.
- [80] P. Newman and K. Ho, "SLAM - Loop closing with visually salient features," in *Robotics and Automation (ICRA 2005), 2005 IEEE International Conference on*, 2005, pp. 635–642.
- [81] A. Gil, O. Reinoso, O. M. Mozos, C. Stachniss, and W. Burgard, "Improving data association in vision-based SLAM," in *Intelligent Robots and Systems, 2006 IEEE/RSJ International Conference on*, 2006, pp. 2076–2081.
- [82] A. Angeli, S. Doncieux, J.-A. Meyer, and D. Filliat, "Incremental vision-based topological slam," in *Intelligent Robots and Systems (IROS 2008), 2008 IEEE/RSJ International Conference on*, 2008, pp. 1031–1036.
- [83] A. Angeli, D. Filliat, S. Doncieux, and J.-A. Meyer, "Fast and incremental method for loop-closure detection using bags of visual words," *IEEE Trans. Robot.*, vol. 24, no. 5, pp. 1027–1037, 2008.
- [84] K. L. Ho and P. Newman, "Detecting loop closure with scene sequences," *Int. J. Comput. Vis.*, vol. 74, no. 3, pp. 261–286, 2007.
- [85] K. Mikolajczyk and C. Schmid, "Indexing based on scale invariant interest points," in *Computer Vision, 2001. ICCV 2001. Proceedings. Eighth IEEE International Conference on*, 2001, vol. 1, pp. 525–531 vol.1.
- [86] A. C. Murillo, J. J. Guerrero, and C. Sagues, "SURF features for efficient robot localization with omnidirectional images," in *Robotics and Automation (ICRA 2007), 2007 IEEE International Conference on*, 2007, pp. 3901–3907.
- [87] M. Cummins and P. Newman, "Appearance-only SLAM at large scale with FAB-MAP 2.0," *Int. J. Rob. Res.*, vol. 30, no. 9, pp. 1100–1123, 2011.
- [88] M. Agrawal, K. Konolige, and M. R. Blas, "CenSurE: Center Surround Extremas for Realtime Feature Detection and Matching," in *Computer Vision - ECCV, 2008*, vol. 4, pp. 102–115.
- [89] C. Mei, G. Sibley, M. Cummins, P. M. Newman, and I. D. Reid, "A constant-time efficient stereo SLAM system," in *BMVC, 2009*, pp. 1–11.

- [90] E. Rosten and T. Drummond, "Machine learning for high-speed corner detection," in *Computer Vision—ECCV 2006*, Springer, 2006, pp. 430–443.
- [91] M. Calonder, V. Lepetit, M. Ozuysal, T. Trzcinski, C. Strecha, and P. Fua, "BRIEF: computing a local binary descriptor very fast," *Pattern Anal. Mach. Intell. IEEE Trans.*, pp. 778–792, 2012.
- [92] J. Sivic and A. Zisserman, "Video Google: A text retrieval approach to object matching in videos," in *IEEE International Conference on Computer Vision (ICCV)*, 2003, pp. 1470–1477.
- [93] L. Fei-Fei and P. Perona, "A Bayesian hierarchical model for learning natural scene categories," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005, vol. 2, pp. 524–531.
- [94] R. Baeza-Yates and B. Ribeiro-Neto, *Modern Information Retrieval*, vol. 463. ACM press New York, 1999.
- [95] D. Nister and H. Stewenius, "Scalable recognition with a vocabulary tree," in *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, 2006, vol. 2, pp. 2161–2168.
- [96] J. Wang, H. Zha, and R. Cipolla, "Combining interest points and edges for content-based image retrieval," in *Image Processing, 2005. ICIP 2005. IEEE International Conference on*, 2005, vol. 3, pp. III–1256–9.
- [97] D. Filliat, "A visual bag of words method for interactive qualitative localization and mapping," in *Robotics and Automation, 2007 IEEE International Conference on*, 2007, pp. 3921–3926.
- [98] R. Paul and P. Newman, "FAB-MAP 3D: Topological Mapping with Spatial and Visual Appearance," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2010, pp. 2649–2656.
- [99] P. Newman, M. Smith, A. Harrison, C. Mei, I. Posner, R. Shade, D. Schröter, W. Churchill, and I. Reid, "Navigating, Recognising and Describing Urban Spaces With Vision and Laser," *Int. J. Rob. Res.*, vol. 28, no. 11–12, pp. 1406–1433, 2009.
- [100] C. Valgren and A. Lilienthal, "SIFT, SURF & seasons: Appearance-based long-term localization in outdoor environments," *Rob. Auton. Syst.*, vol. 58, no. 2, pp. 157–165, 2010.
- [101] C. Cadena and J. Neira, "A learning algorithm for place recognition," *ICRA 2011 Workshop on Long-term Autonomy*. Shanghai, China, 2011.
- [102] E. Johns and G.-Z. Yang, "Feature co-occurrence maps: Appearance-based localisation throughout the day," in *Proc. ICRA*, 2013.
- [103] M. Milford, W. Scheirer, E. Vig, A. Glover, O. Baumann, J. Mattingley, and D. Cox, "Condition-invariant, top-down visual place recognition," in *Robotics and Automation (ICRA 2014), 2014 IEEE International Conference on*, 2014.
- [104] B. J. A. Kröse, N. Vlassis, R. Bunschoten, and Y. Motomura, "A Probabilistic Model for Appearance-Based Robot Localization," in *In First European Symposium on Ambient Intelligence (EUSAI, 2000*, pp. 264–274.
- [105] P. Lamon, I. Nourbakhsh, B. Jensen, and R. Siegwart, "Deriving and matching image fingerprint sequences for mobile robot localization," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2001, pp. 1609–1614.
- [106] J. Canny, "A computational approach to edge detection," *Pattern Anal. Mach. Intell. IEEE Trans.*, no. 6, pp. 679–698, 1986.
- [107] C. Harris and M. Stephens, "A Combined Corner and Edge Detector," in *4th Alvey Vision Conference*, 1988, vol. 147–151.
- [108] H. Badino, D. Huber, and T. Kanade, "Real-time topometric localization," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2012, pp. 1635–1642.
- [109] N. Sunderhauf and P. Protzel, "BRIEF-Gist - closing the loop by simple means," in *Intelligent Robots and Systems (IROS 2011), 2011 IEEE/RSJ International Conference on*, 2011, pp. 1234–1241.
- [110] A. Murillo and J. Košecká, "Experiments in place recognition using gist panoramas," in *IEEE International Conference on Computer Vision (ICCV)*, 2009, pp. 2196–2203.
- [111] C. Siagian and L. Itti, "Biologically inspired mobile robot vision localization," *Robot. IEEE Trans.*, vol. 25, no. 4, pp. 861–873, 2009.
- [112] G. Singh and J. Košecká, "Visual loop closing using gist descriptors in manhattan world," in *Omnidirectional Robot Vision workshop, with IEEE Int. Conf. on Robotics and Automation*, 2010.
- [113] Y. Liu and H. Zhang, "Visual loop closure detection with a compact image descriptor," in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, 2012, pp. 1051–1056.
- [114] C. Mei, G. Sibley, and P. Newman, "Closing loops without places," in *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*, 2010, pp. 3738–3744.
- [115] S. Lynen, M. Bosse, P. Furgale, and R. Siegwart, "Placeless Place-Recognition," in *International Conference on 3D Vision*, 2014.
- [116] M. Milford and G. Wyeth, "Mapping a Suburb With a Single Camera Using a Biologically Inspired SLAM System," *Robot. IEEE Trans.*, vol. 24, no. 5, pp. 1038–1053, 2008.
- [117] P. Furgale and T. Barfoot, "Visual teach and repeat for long-range rover autonomy," *J. F. Robot.*, vol. 27, no. 5, pp. 534–560, 2010.
- [118] M. Milford and G. Wyeth, "SeqSLAM: Visual route-based navigation for sunny summer days and stormy winter nights," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2012, pp. 1643–1649.
- [119] T. Naseer, L. Spinello, W. Burgard, and C. Stachniss, "Robust visual robot localization across seasons using network flows," in *Conference on the Association for the Advancement of Artificial Intelligence*, 2014.
- [120] C. McManus, B. Upcroft, and P. Newman, "Scene signatures: Localised and point-less features for localisation," in *Robotics: Science and Systems*, 2014.
- [121] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, 2005, vol. 1, pp. 886–893.
- [122] 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," in *Robotics: Science and Systems*, 2015.
- [123] C. L. Zitnick and P. Dollár, "Edge Boxes: Locating Object Proposals from Edges," in *ECCV*, 2014.
- [124] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *NIPS*, 2012, vol. 1, no. 2, p. 4.
- [125] B. Krose, N. Vlassis, R. Bunschoten, and Y. Motomura, "A probabilistic model for appearance-based robot localization," *Image Vis. Comput.*, vol. 19, no. 6, pp. 381–391, 2001.
- [126] J. Neira, A. J. Davison, and J. J. Leonard, "Guest editorial - Special issue on visual SLAM," *IEEE Trans. Robot.*, vol. 24, no. 5, pp. 929–931, 2008.
- [127] F. Endres, J. Hess, N. Engelhard, J. Sturm, D. Cremers, and W. Burgard, "An evaluation of the RGB-D SLAM system," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, 2012, pp. 1691–1696.
- [128] M. Labbe and F. Michaud, "Online global loop closure detection for large-scale multi-session graph-based SLAM," in *Intelligent Robots and Systems (IROS 2014), 2014 IEEE/RSJ International Conference on*, 2014, pp. 2661–2666.
- [129] T. Whelan, M. Kaess, R. Finman, M. Fallon, H. Johannsson, J. Leonard, and J. McDonald, "Real-time large scale dense RGB-D SLAM with volumetric fusion," *Int. J. Rob. Res.*, 2014.
- [130] A. J. Davison and D. W. Murray, "Simultaneous localization and map-building using active vision," *Pattern Anal. Mach. Intell. IEEE Trans.*, vol. 24, no. 7, pp. 865–880, 2002.
- [131] K. Konolige and J. Bowman, "Towards lifelong visual maps," in *Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on*, 2009, pp. 1156–1163.
- [132] K. Konolige, J. Bowman, J. D. Chen, P. Mihelich, M. Calonder, V. Lepetit, and P. Fua, "View-based maps," *Int. J. Rob. Res.*, vol. 29, no. 8, pp. 941–957, 2010.
- [133] B. Triggs, P. F. McLauchlan, R. I. Hartley, and A. W. Fitzgibbon, "Bundle adjustment—a modern synthesis," in *Vision algorithms: theory and practice*, Springer, 2000, pp. 298–372.
- [134] G. Klein and D. Murray, "Parallel tracking and mapping for small AR workspaces," in *Mixed and Augmented Reality, 2007. ISMAR 2007. 6th IEEE and ACM International Symposium on*, 2007, pp. 225–234.
- [135] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, "DTAM: Dense tracking and mapping in real-time," in *Computer Vision (ICCV), 2011 IEEE International Conference on*, 2011, pp. 2320–2327.

- [136] J. Engel, T. Schops, J. Sturm, and D. Cremers, "LSD-SLAM: Large-Scale Direct Monocular SLAM," in *European Conference on Computer Vision (ECCV)*, 2014.
- [137] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardós, "ORB-SLAM: a versatile and accurate monocular SLAM system," *CoRR*, vol. abs/1502.0, 2015.
- [138] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "ORB: An efficient alternative to SIFT or SURF," in *Proceedings of the IEEE International Conference on Computer Vision*, 2011, pp. 2564–2571.
- [139] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. W. Fitzgibbon, "KinectFusion: Real-time dense surface mapping and tracking," in *10th IEEE International Symposium on Mixed and Augmented Reality, ISMAR 2011, Basel, Switzerland, October 26-29, 2011*, 2011, pp. 127–136.
- [140] R. Finman, L. Paull, and J. J. Leonard, "Toward object-based place recognition in dense RGB-D maps," in *ICRA workshop on visual place recognition in changing environments*, 2015.
- [141] M. Bosse, P. Newman, J. Leonard, and S. Teller, "Simultaneous localization and map building in large-scale cyclic environments using the Atlas framework," *Int. J. Rob. Res.*, vol. 23, no. 12, pp. 1113–1139, 2004.
- [142] F. Li and J. Koščeká, "Probabilistic location recognition using reduced feature set," in *Robotics and Automation (ICRA 2006), 2006 IEEE International Conference on*, 2006, pp. 3405–3410.
- [143] M. Cummins and P. Newman, "Highly scalable appearance-only SLAM - FAB-MAP 2.0," *Robot. Sci. Syst.*, vol. 1, pp. 12–18, 2009.
- [144] M. Mohan, D. Gálvez-López, C. Monteleoni, and G. Sibley, "Environment selection and hierarchical place recognition," in *Robotics and Automation (ICRA 2015), 2015 IEEE International Conference on*, 2015.
- [145] H. Andreasson, A. Treptow, and T. Duckett, "Localization for Mobile Robots using Panoramic Vision, Local Features and Particle Filter," *Proc. 2005 IEEE Int. Conf. Robot. Autom.*, 2005.
- [146] W. Maddern, M. Milford, and G. Wyeth, "Capping computation time and storage requirements for appearance-based localization with CAT-SLAM," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, 2012, pp. 822–827.
- [147] D. L. Donoho and Y. Tsaig, "Fast solution of-norm minimization problems when the solution may be sparse," *Inf. Theory, IEEE Trans.*, vol. 54, no. 11, pp. 4789–4812, 2008.
- [148] M. Milford, "Visual route recognition with a handful of bits," *Proc. Robot. Sci. Syst. Conf. 2012*, 2012.
- [149] M. Milford, "Vision-based place recognition: How low can you go?," *Int. J. Rob. Res.*, vol. 32, no. 7, pp. 766–789, 2013.
- [150] E. Pepperell, P. Corke, and M. Milford, "All-environment visual place recognition with SMART," in *Robotics and Automation (ICRA 2014), 2014 IEEE International Conference on*, 2014, pp. 1612–1618.
- [151] T. Duckett, S. Marsland, and J. Shapiro, "Learning globally consistent maps by relaxation," in *Robotics and Automation, 2000. Proceedings. ICRA'00. IEEE International Conference on*, 2000, vol. 4, pp. 3841–3846.
- [152] D. Filliat and J.-A. Meyer, "Global localization and topological map-learning for robot navigation," *From Anim. to Animat.*, vol. 7, pp. 131–140, 2002.
- [153] A. Angeli, S. Doncieux, J.-A. Meyer, and D. Filliat, "Visual topological SLAM and global localization," in *Robotics and Automation, 2009. ICRA'09. IEEE International Conference on*, 2009, pp. 4300–4305.
- [154] S. Bazeille and D. Filliat, "Incremental topo-metric SLAM using vision and robot odometry," in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, 2011, pp. 4067–4073.
- [155] C. Mei, G. Sibley, M. Cummins, P. Newman, and I. Reid, "RSLAM: A system for large-scale mapping in constant-time using stereo," *Int. J. Comput. Vis.*, vol. 94, no. 2, pp. 198–214, 2011.
- [156] G. Sibley, C. Mei, I. Reid, and P. Newman, "Adaptive relative bundle adjustment," in *Robotics: Science and Systems*, 2009.
- [157] K. Konolige, E. Marder-Eppstein, and B. Marthi, "Navigation in hybrid metric-topological maps," in *Robotics and Automation (ICRA 2011), 2011 IEEE International Conference on*, 2011, pp. 3041–3047.
- [158] P. Beeson, J. Modayil, and B. Kuipers, "Factoring the mapping problem: Mobile robot map-building in the Hybrid Spatial Semantic Hierarchy," *Int. J. Rob. Res.*, vol. 29, no. 4, pp. 428–459, 2010.
- [159] H. Johannsson, M. Kaess, M. F. Fallon, and J. J. Leonard, "Temporally scalable visual SLAM using a reduced pose graph," in *Robotics and Automation (ICRA 2013), 2013 IEEE International Conference on*, 2013.
- [160] I. Mahon, S. Williams, O. Pizarro, and M. Johnson-Roberson, "Efficient view-based SLAM using visual loop closures," *Robot. IEEE Trans.*, vol. 24, no. 5, pp. 1002–1014, 2008.
- [161] R. M. Eustice, H. Singh, J. J. Leonard, and M. R. Walter, "Visually mapping the RMS Titanic: Conservative covariance estimates for SLAM information filters," *Int. J. Rob. Res.*, vol. 25, no. 12, pp. 1223–1242, Dec. 2006.
- [162] F. Lu and E. Milios, "Globally consistent range scan alignment for environment mapping," *Auton. Robots*, vol. 4, no. 4, pp. 333–349, 1997.
- [163] J. S. Gutmann and K. Konolige, "Incremental mapping of large cyclic environments," in *IEEE International Symposium on Computational Intelligence in Robotics and Automation*, 1999, pp. 318–325.
- [164] H. Durrant-Whyte and T. Bailey, "Simultaneous Localization and Mapping: Part I," *Robot. Autom. Mag. IEEE*, vol. 13, no. 2, pp. 99–110, 2006.
- [165] T. Bailey and H. Durrant-Whyte, "Simultaneous Localization and Mapping (SLAM): Part II," *Robot. Autom. Mag. IEEE*, vol. 13, no. 3, pp. 108–117, 2006.
- [166] S. Thrun and J. Leonard, "Simultaneous Localisation and Mapping," in *Springer Handbook of Robotics*, B. Siciliano and O. Khatib, Eds. 2008, pp. 871–889.
- [167] R. C. Smith and P. Cheeseman, "On the representation and estimation of spatial uncertainty," *Int. J. Rob. Res.*, vol. 5, no. 4, pp. 56–68, 1987.
- [168] J. E. Guivant and E. M. Nebot, "Optimization of the simultaneous localization and map-building algorithm for real-time implementation," *IEEE Trans. Robot. Autom.*, vol. 17, no. 3, pp. 242–257, 2001.
- [169] M. W. M. G. Dissanayake, P. Newman, S. Clark, H. F. Durrant-Whyte, and M. Csorba, "A solution to the simultaneous localization and map building (SLAM) problem," *IEEE Trans. Robot. Autom.*, vol. 17, no. 3, pp. 229–241, 2001.
- [170] M. Montemerlo, S. Thrun, and B. Siciliano, *FastSLAM: a scalable method for the simultaneous localization and mapping problem in robotics*, no. v 27. Berlin: Springer, 2007.
- [171] S. Thrun and M. Montemerlo, "The GraphSLAM algorithm with applications to large-scale mapping of urban structures," *Int. J. Rob. Res.*, vol. 25, no. 5–6, pp. 403–429, 2006.
- [172] E. Olson, J. Leonard, and S. Teller, "Fast iterative alignment of pose graphs with poor initial estimates," in *Robotics and Automation (ICRA 2006), 2006 IEEE International Conference on*, 2006, pp. 2262–2269.
- [173] M. Kaess, A. Ranganathan, and F. Dellaert, "iSAM: Incremental smoothing and mapping," *IEEE Trans. Robot.*, vol. 24, no. 6, pp. 1365–1378, 2008.
- [174] E. Eade and T. Drummond, "Edge landmarks in monocular SLAM," in *BMVC*, 2006, pp. 7–16.
- [175] E. Eade and T. Drummond, "Scalable monocular SLAM," in *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, 2006, vol. 1, pp. 469–476.
- [176] E. Eade and T. Drummond, "Unified loop closing and recovery for real time monocular SLAM," *BMVC*, vol. 13, p. 136, 2008.
- [177] A. Jacobson, Z. Chen, and M. Milford, "Autonomous Multisensor Calibration and Closed-loop Fusion for SLAM," *J. F. Robot.*, 2014.
- [178] D. Galvez-Lopez and J. D. Tardos, "Real-time loop detection with bags of binary words," in *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on*, 2011, pp. 51–58.
- [179] N. Vlassis and B. Krose, "Robot environment modeling via principal component regression," in *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, 1999, pp. 677–682.
- [180] N. Vlassis and B. Krose, "Mixture conditional density estimation with the EM algorithm," in *ICANN'99: 9th International Conference on Artificial Neural Networks*, 1999.
- [181] B. D. Ripley, *Pattern Recognition and Neural Networks*. Cambridge University Press, 1996.

- [182] A. Angeli, S. Doncieux, J. A. Meyer, and D. Filliat, "Real-time visual loop-closure detection," in *Robotics and Automation (ICRA 2008)*, 2008 IEEE International Conference on, 2008, pp. 1842–1847.
- [183] C. Siagian and L. Itti, "Biologically-inspired robotics vision monte-carlo localization in the outdoor environment," in *Intelligent Robots and Systems (IROS 2007)*, 2007 IEEE/RSJ International Conference on, 2007, pp. 1723–1730.
- [184] E. Garcia-Fidalgo and A. Ortiz, "Probabilistic appearance-based mapping and localization using visual features," in *IbPRIA 2013: Iberian Conference on Pattern Recognition and Image Analysis*, 2013, pp. 277–285.
- [185] C. K. Chow and C. N. Liu, "Approximating discrete probability distributions with dependence trees," *IEEE Trans. Inf. Theory*, vol. 14, no. 3, pp. 462–467, 1968.
- [186] R. Paul and P. Newman, "Self-help: Seeking out perplexing images for ever improving topological mapping," *Int. J. Rob. Res.*, vol. 32, no. 14, pp. 1742–1766, 2013.
- [187] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, 2003.
- [188] E. Olson, "Recognizing places using spectrally clustered local matches," *Rob. Auton. Syst.*, vol. 57, no. 12, pp. 1157–1172, Dec. 2009.
- [189] C. Giovannangeli, P. Gaussier, and G. Desilles, "Robust mapless outdoor vision-based navigation," in *Intelligent Robots and Systems (IROS 2006)*, IEEE/RSJ International Conference on, 2006, pp. 3293–3300.
- [190] E. Olson and P. Agarwal, "Inference on networks of mixtures for robust robot mapping," *Robotics: Science and Systems*. Sydney, Australia, 2012.
- [191] N. Sunderhauf and P. Protzel, "Switchable constraints for robust pose graph SLAM," in *Intelligent Robots and Systems (IROS 2012)*, 2012 IEEE/RSJ International Conference on, 2012, pp. 1879–1884.
- [192] Y. Latif, C. Cadena, and J. Neira, "Robust loop closing over time for pose graph SLAM," *Int. J. Rob. Res.*, vol. 32, no. 14, pp. 1611–1626, 2013.
- [193] C. McManus, W. Churchill, W. Maddern, A. Stewart, and P. Newman, "Shady dealings: robust, long-term visual localisation using illumination invariance," in *Robotics and Automation (ICRA 2014)*, 2014 IEEE International Conference on, 2014.
- [194] T. Barfoot, J. Kelly, and G. Sibley, "Special Issue on Long-Term Autonomy," *Int. J. Rob. Res.*, vol. 32, no. 14, pp. 1609–1610, 2013.
- [195] C. Valgren and A. Lilienthal, "SIFT, SURF and Seasons: Long-term Outdoor Localization Using Local Features," in *European Conference on Mobile Robotics (ECMR)*, 2007, vol. 128, pp. 1–6.
- [196] P. Ross, A. English, D. Ball, B. Upcroft, G. Wyeth, and P. Corke, "A novel method for analysing lighting variance," in *Proceedings of Australasian Conference on Robotics and Automation*, 2013, pp. 1–8.
- [197] P. Ross, A. English, D. Ball, and P. Corke, "A method to quantify a descriptor's illumination variance," in *Proceedings of Australasian Conference on Robotics and Automation*, 2014.
- [198] M. J. Milford, I. Turner, and P. Corke, "Long exposure localization in darkness using consumer cameras," in *Robotics and Automation (ICRA 2013)*, 2013 IEEE International Conference on, 2013, pp. 3755–3761.
- [199] N. Sünderhauf, P. Neubert, and P. Protzel, "Are we there yet? Challenging SeqSLAM on a 3000 km journey across all four seasons," in *Proc. of Workshop on Long-Term Autonomy*, IEEE International Conference on Robotics and Automation (ICRA), 2013.
- [200] S. Nuske, J. Roberts, and G. Wyeth, "Robust outdoor visual localization using a three-dimensional-edge map," *J. F. Robot.*, vol. 26, no. 9, pp. 728–756, 2009.
- [201] P. Borges, R. Zlot, M. Bosse, S. Nuske, and A. Tews, "Vision-based localization using an edge map extracted from 3D laser range data," in *Robotics and Automation (ICRA)*, 2010 IEEE International Conference on, 2010, pp. 4902–4909.
- [202] P. Corke, R. Paul, W. Churchill, and P. Newman, "Dealing with shadows: Capturing intrinsic scene appearance for image-based outdoor localisation," in *Intelligent Robots and Systems (IROS 2013)*, 2013 IEEE/RSJ International Conference on, 2013, pp. 2085–2092.
- [203] C. McManus, P. Furgale, and T. D. Barfoot, "Towards lighting-invariant visual navigation: An appearance-based approach using scanning laser-range finders," *Rob. Auton. Syst.*, 2013.
- [204] W. Maddern and S. Vidas, "Towards robust night and day place recognition using visible and thermal imaging," *RSS 2012 Beyond laser Vis. Altern. Sens. Tech. Robot. Percept.*, 2012.
- [205] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN Features off-the-shelf: an Astounding Baseline for Recognition," in *Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2014 IEEE Conference on, 2014, pp. 512–519.
- [206] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Learning and transferring mid-level image representations using convolutional neural networks," in *Computer Vision and Pattern Recognition (CVPR)*, 2014 IEEE Conference on, 2014, pp. 1717–1724.
- [207] L. O. J. A. & M. M. Chen Zetao, "Convolutional Neural Network-based Place Recognition," in *Australasian Conference on Robotics and Automation 2014*, 2014.
- [208] N. Sünderhauf, S. Shirazi, F. Dayoub, B. Upcroft, and M. Milford, "On the performance of ConvNet features for place recognition," in *Intelligent Robots and Systems (IROS 2015)*, 2015 IEEE/RSJ International Conference on, 2015.
- [209] W. Maddern, A. D. Stewart, C. McManus, B. Upcroft, W. Churchill, and P. Newman, "Illumination invariant imaging: applications in robust vision-based localisation, mapping and classification for autonomous vehicles," *Proc. Work. Vis. Place Recognit. Chang. Environ. IEEE Int. Conf. Robot. Autom.*, 2014.
- [210] N. Jacobs, N. Roman, and R. Pless, "Consistent Temporal Variations in Many Outdoor Scenes," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Minneapolis, MN, 2007.
- [211] A. Abrams, E. Feder, and R. Pless, "Exploratory analysis of time-lapse imagery with fast subset PCA," in *Applications of Computer Vision (WACV)*, 2011 IEEE Workshop on, 2011, pp. 336–343.
- [212] A. Ranganathan, S. Matsumoto, and D. Ilstrup, "Towards illumination invariance for visual localization," in *Robotics and Automation (ICRA 2013)*, 2013 IEEE International Conference on, 2013, pp. 3791–3798.
- [213] A. Mikulík, M. Perdoch, O. Chum, and J. Matas, "Learning a fine vocabulary," in *Computer Vision—ECCV 2010*, Springer, 2010, pp. 1–14.
- [214] N. Carlevaris-Bianco and R. M. Eustice, "Learning visual feature descriptors for dynamic lighting conditions," *Proc. of Workshop on Visual Place Recognition in Changing Environments*, IEEE International Conference on Robotics and Automation (ICRA). 2014.
- [215] J. M. Bromley, I. Guyon, Y. LeCun, E. Sackinger, and R. Shah, "Signature verification using a Siamese time delay neural network," in *7th Annual Neural Information Processing Systems Conference*, 1994, vol. 7, pp. 737–744.
- [216] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Susstrunk, "SLIC superpixels Compared to state-of-the-art superpixel methods," *Pattern Anal. Mach. Intell. IEEE Trans.*, vol. 34, no. 11, pp. 2274–2282, 2012.
- [217] S. Lowry, M. Milford, and G. Wyeth, "Transforming morning to afternoon using linear regression techniques," in *Robotics and Automation (ICRA 2014)*, 2014 IEEE International Conference on, 2014.
- [218] S. Lowry, G. Wyeth, and M. Milford, "Unsupervised online learning of condition-invariant images for place recognition," *Proc. Australas. Conf. Robot. Autom. 2014*, 2014.
- [219] P. Biber and T. Duckett, "Dynamic maps for long-term operation of mobile service robots," in *Robotics: Science and Systems*, 2005, pp. 17–24.
- [220] F. Dayoub and T. Duckett, "An adaptive appearance-based map for long-term topological localization of mobile robots," in *Intelligent Robots and Systems (IROS 2008)*, 2008 IEEE/RSJ International Conference on, 2008, pp. 3364–3369.
- [221] T. Morris, F. Dayoub, P. Corke, G. Wyeth, and B. Upcroft, "Multiple map hypotheses for planning and navigating in non-stationary environments," in *Robotics and Automation (ICRA 2014)*, 2014 IEEE International Conference on, 2014.
- [222] T. Bailey, "Mobile robot localisation and mapping in extensive outdoor environments," University of Sydney, 2002.

- [223] J. Andrade-Cetto and A. Sanfeliu, "Concurrent map building and localisation in indoor dynamic environments," *Int. J. Pattern Recognit. Artif. Intell.*, vol. 16, no. 3, pp. 361–374, 2002.
- [224] A. H. Hafez, M. Singh, K. M. Krishna, and C. V. Jawahar, "Visual localization in highly crowded urban environments," in *Intelligent Robots and Systems (IROS 2013), 2013 IEEE/RSJ International Conference on*, 2013, pp. 2778–2783.
- [225] E. Johns and G.-Z. Yang, "Generative methods for long-term place recognition in dynamic scenes," *Int. J. Comput. Vis.*, vol. 106, no. 3, pp. 297–314, 2014.
- [226] M. Milford and G. Wyeth, "Persistent navigation and mapping using a biologically inspired SLAM system," *Int. J. Rob. Res.*, vol. 29, no. 9, pp. 1131–1153, 2010.
- [227] P. Biber and T. Duckett, "Experimental analysis of sample-based maps for long-term SLAM," *Int. J. Rob. Res.*, vol. 28, no. 1, pp. 20–33, 2009.
- [228] C. Stachniss and W. Burgard, "Mobile robot mapping and localization in non-static environments," in *AAAI'05 - The 20th National Conference on Artificial Intelligence*, 2005, pp. 1324–1329.
- [229] W. Churchill and P. Newman, "Practice makes perfect? Managing and leveraging visual experiences for lifelong navigation," in *Robotics and Automation (ICRA 2012), 2012 IEEE International Conference on*, 2012, pp. 4525–4532.
- [230] C.-C. Wang, C. Thorpe, S. Thrun, M. Hebert, and H. Durrant-Whyte, "Simultaneous localization, mapping and moving object tracking," *Int. J. Rob. Res.*, vol. 26, no. 9, pp. 889–916, 2007.
- [231] J. F. Dong, S. Wijesoma, and A. P. Shacklock, "Extended Rao-Blackwellised genetic algorithmic filter SLAM in dynamic environment with raw sensor measurement," in *Intelligent Robots and Systems (IROS 2007), 2007 IEEE/RSJ International Conference on*, 2007, pp. 1473–1478.
- [232] D. F. Wolf and G. S. Sukhatme, "Mobile robot simultaneous localization and mapping in dynamic environments," *Auton. Robots*, vol. 19, no. 1, pp. 53–65, 2005.
- [233] D. Meyer-Delius, J. Hess, G. Grisetti, and W. Burgard, "Temporary Maps for Robust Localisation in Semi-static environments," *Intelligent Robots and Systems (IROS 2010), IEEE/RSJ International Conference on*, 2010.
- [234] T. Krajník, J. P. Fentanes, O. M. Mozos, T. Duckett, J. Ekekrantz, and M. Hanheide, "Long-term topological localisation for service robots in dynamic environments using spectral maps," in *Intelligent Robots and Systems (IROS 2014), 2014 IEEE/RSJ International Conference on*, 2014, pp. 4537–4542.
- [235] N. Carlevaris-Bianco and R. M. Eustice, "Learning temporal co-observability relationships for lifelong robotic mapping," in *IROS Workshop on Lifelong Learning for Mobile Robotics Applications*, 2012.
- [236] C. Linegar, W. Churchill, and P. Newman, "Work smart, not hard: Recalling relevant experiences for vast-scale but time-constrained localisation," in *Robotics and Automation (ICRA 2015), 2015 IEEE International Conference on*, 2015, pp. 90–97.
- [237] J. Lafferty, A. McCallum, and F. C. N. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data," 2001.
- [238] Y. Liu and H. Zhang, "Towards improving the efficiency of sequence-based SLAM," in *Mechatronics and Automation (ICMA 2013), 2013 IEEE International Conference on*, 2013, pp. 1261–1266.
- [239] P. Hansen and B. Browning, "Visual place recognition using HMM sequence matching," in *Intelligent Robots and Systems (IROS 2014), 2014 IEEE/RSJ International Conference on*, 2014, pp. 4549–4555.



Stephanie Lowry received the B.Sc (Hons) degree in logic and computation and the M.Sc in computer science, both from Victoria University of Wellington, New Zealand, and the Ph.D in engineering from Queensland University of Technology, Australia in 2014.

She is currently a Postdoctoral Researcher with the Australian Research

Council Centre of Excellence in Robotic Vision, Queensland University of Technology. Her research interests include long-term robot autonomy and the application of machine learning to robotic vision.



Niko Sünderhauf received the Ph.D. in 2012 and the Diploma (M.Sc.) in computer science in 2006, both from Technische Universität Chemnitz, Germany, where he was a research fellow between 2006 and 2014.

In March 2014 he joined Queensland University of Technology in Brisbane and is now a research fellow with the Australian Research Council Centre of Excellence in Robotic Vision. His research interests include robust visual perception, place recognition in changing conditions, deep learning, SLAM, long-term autonomy, and probabilistic estimation with graphical models. Apart from mobile robotics, his research covers robust methods for sensor fusion, especially NLOS-mitigation for satellite-based localization systems.



Paul Newman (F'15) received the M.Eng. degree in engineering science from the University of Oxford, Oxford, U.K., in 1995 and the Ph.D. degree in autonomous navigation from the Australian Center for Field Robotics, University of Sydney, Sydney, N.S.W., Australia.

He is currently a BP Professor of information engineering with the Department of Engineering Science, University of Oxford. He heads the Oxford Mobile Robotics Research Group. His research interest includes autonomous navigation, especially over large spatial and temporal scales. In 1999, he returned to the U.K. to work in the commercial sub-sea navigation industry. In late 2000, he joined the Department of Ocean Engineering, Massachusetts Institute of Technology, Cambridge, where as a Postdoctoral Researcher and later a Research Scientist, he worked on algorithms and software for robust autonomous navigation for both land and sub-sea agents. In early 2003, he returned to the University of Oxford as a Departmental Lecturer in engineering science before being appointed to a University Lectureship in 2005.



John J. Leonard (S'87-M'87-F'13) received the B.S.E. degree in electrical engineering and science from the University of Pennsylvania, Philadelphia, and the D.Phil. degree in engineering science from the University of Oxford, Oxford, U.K., in 1994.

Currently, he is a Professor of Mechanical and Ocean Engineering in the Department of Mechanical Engineering, Massachusetts Institute of Technology (MIT), Cambridge. He is also a member of the MIT Computer Science and Artificial Intelligence Laboratory (CSAIL). His research addresses the problems of navigation, mapping, and persistent autonomy for

autonomous mobile robots operating in unstructured environments.



David Cox received the A.B. degree in biology and psychology from Harvard University, Cambridge, MA, USA, in 2000 and the Ph.D. degree in computational neuroscience from Massachusetts Institute of Technology, Cambridge, MA, USA in 2007.

He is currently an Assistant Professor of Molecular and Cellular Biology and of Computer Science at Harvard University. He was previously a Junior Fellow at the Rowland Institute at Harvard University, a multidisciplinary institute focused on high-risk, high-reward scientific research at the boundaries of traditional fields. His laboratory seeks to understand the computational underpinnings of high-level visual processing through concerted efforts in both reverse- and forward-engineering. To this end, his group employs a wide range of experimental techniques (ranging from microelectrode recordings in living brains to visual psychophysics in humans) to probe natural systems, while at the same time actively developing practical computer vision systems based on what is learned about the brain.



Peter Corke received Bachelor of Engineering, Masters of Engineering and PhD degrees from University of Melbourne, Melbourne, Australia and is a Fellow of the IEEE.

He is currently director of the Australian Research Council Centre of Excellence in Robotic Vision, and a Professor of Robotics and Control at Queensland University of Technology (QUT), Brisbane, Australia. His research is concerned with robotic vision, flying robots and robots for agriculture.

Prof. Corke worked at the University of Melbourne, first as a research assistant and later as a lecturer. In 1984 he commenced with the Commonwealth Scientific and Industrial Research Organisation (CSIRO), where he founded the Autonomous Systems laboratory of the CSIRO ICT Centre, and served as Research Director from 2004–2007. He was a Senior Principal Research Scientist when he left to take up a chair at QUT in 2010. He was editor-in-chief of the IEEE Robotics & Automation magazine from 2009 to 2013, and is a founding editor of the Journal of Field Robotics.



Michael Milford received the Ph.D. in electrical engineering and the Bachelor of Mechanical and Space Engineering from the University of Queensland (UQ), Brisbane, Australia.

He is currently a Senior Lecturer and Australian Research Council Future Fellow at Queensland University of Technology, Brisbane, Australia, and a Chief Investigator for the Australian Centre of Excellence for Robotic Vision. He was a Research Fellow on the Thinking Systems Project at the

Queensland Brain Institute on the Thinking Systems Project until 2010, when he became a Lecturer at QUT.

Dr. Milford was awarded an inaugural Australian Research Council Discovery Early Career Researcher Award in 2012 and became a Microsoft Faculty Fellow in 2013.