# Hybrid Belief Pruning with Guarantees for Viewpoint-Dependent Semantic SLAM

Tuvy Lemberg and Vadim Indelman

*Abstract*— Semantic simultaneous localization and mapping is a subject of increasing interest in robotics and AI that directly influences the autonomous vehicles industry, the army industries, and more. One of the challenges in this field is to obtain object classification jointly with robot trajectory estimation. Considering view-dependent semantic measurements, there is a coupling between different classes, resulting in a combinatorial number of hypotheses. A common solution is to prune hypotheses that have a sufficiently low probability and to retain only a limited number of hypotheses. However, after pruning and renormalization, the updated probability is overconfident with respect to the original probability. This is especially problematic for systems that require high accuracy. If the prior probability of the classes is independent, the original normalization factor can be computed efficiently without pruning hypotheses. To the best of our knowledge, this is the first work to present these results. If the prior probability of the classes is dependent, we propose a lower bound on the normalization factor that ensures cautious results. The bound is calculated incrementally and with similar efficiency as in the independent case. After pruning and updating based on the bound, this belief is shown empirically to be close to the original belief.

## I. Introduction

In robotics, the problem of simultaneous localization and mapping (SLAM) concerned with estimating robot poses and, concurrently, reconstructing the environment observed thus far by the robot [1], [2], [3], [4]. The SLAM problem is essential for the robot to be able to navigate autonomously in uncertain or unknown environments.

Recent advances in object recognition and classification enabled to incorporate object classification and semantic features within the SLAM framework, in order to improve both localization estimations and classification [5], [6], [7]. In particular, semantic features are considered to be more robust and indicative than geometric features and can be used in ambiguous and repetitive environments, see e.g. [8]. In addition, semantic SLAM provides a high level understanding that is essential for autonomous decision making and advanced robotic capabilities.

Many semantic SLAM works make use of only the maximum likelihood of a classifier output [9], [10], [11]. This approach reduces the classifier output to one discrete variable and cannot represent the uncertainty of the classifier given a specific image input. More sophisticated approaches consider that classes are not fixed to the maximum likelihood.

Yet, many of these approaches assume that object's class and object's poses are independent, and as such, solve two separable problems (e.g. [5], [6]). However, this corresponds to not taking into account the inherent coupling between classifier output and the viewpoint from which the object is seen.

In recent years, several works considered this coupling by using a viewpoint-dependent semantic observation model. In particular, it has been shown that utilizing such a model enhances disambiguation and localization in challenging scenarios [7]. The resulting belief is hybrid, combining continuous variables such as robot and object poses with discrete variables such as object classes. Moreover, the statistical dependency between the classifier output and the viewpoint creates a statistical dependency between classes of different objects. Another source for this coupling is a dependent prior probability, i.e. when the prior over classification variables of all objects is *not* equal to the product of priors on each of the objects. Such a case is more general as it allows to incorporate readily available statistical knowledge (e.g. high chances to find a keyboard next to a monitor). Therefore, the number of hypotheses is the number of class combinations, which increases *exponentially* with the number of objects. With a few dozen objects and a few hundred classes, the computational complexity becomes prohibitively expensive.

A common way to handle this prohibitive computational complexity is by pruning hypotheses that are considered with the lowest probability. Yet, after pruning and renormalization, the resulting belief is overconfident in its hypotheses with respect to the original belief. Overconfident probabilities may lead to reckless and dangerous behavior. An autonomous vehicle that considers to continue driving according to the overconfident classification may cause fatal errors. Additionally, there is no indication of the significance of pruned hypotheses after they have been pruned. Yet, such an indication can be of prime importance for safe autonomy and the decision making process; furthermore, it may be used to trigger a change in maintained hypotheses (e.g. resurrection of a previously pruned hypothesis).

In this paper, we suggest two methods for maintaining probabilities after pruning in a more realistic and conservative manner. We show that the in case of an a-priori independent class probability, the original normalization factor without pruning can be calculated very efficiently. By obtaining the original normalization factor, we can retrieve the *exact* probability of each hypothesis separately. If the a-priori probability is dependent, we propose a lower bound on the probability of the retained hypotheses. Applying both

methods, provides an indication of the probability of the pruned hypotheses. Using the first method, the probability of the entire pruned set can be computed. Using the second method, the probability of the entire pruned set can be bounded from above. In terms of computational complexity, both methods are similar to pruning.

## II. RELATED WORK

It appears that most semantic SLAM methods can be divided into two types. In the first type, the semantic observations are considered class-dependent and viewpoint-independent. In the second type, the semantic observations are considered to be class independent and viewpoint dependent. There are very few studies which consider the semantic observations to be both class-dependent and viewpoint-dependent. Part of the reason for this is the high computational complexity involved in using this model.

Omidshafiei et al. [12] proposed a Dirichlet distribution as the noisy semantic observation model for the classifier output. They developed a Bayesian filter and showed that the resulting classification is more robust to ambiguity. The semantic observations are assumed to be viewpoint-independent. Nicholson et al. [13] used two-dimensional object detection and represented objects as three-dimensional oriented quadratic surfaces. The semantic observations are the object shape detections which are considered to be viewpoint-dependent and class-independent; Classification took place externally. Yang and Scherer, [10] proposed object SLAM fed by a single image 3D cuboid object detector - the semantic model, which is viewpoint-dependent and class-independent. In Doherty et al. [5] for each new measurement, the marginal posterior of the data association was computed. The posterior of the data association was then used to compute the posterior of the remaining parameters. As a result, they did not have to directly maintain a hybrid belief. Moreover, they implemented this method within a factor graph in [14]. In both works, the semantic model is assumed to be view independent, thereby avoiding the mentioned high computational complexity.

Both viewpoint-dependent and class-dependent models were recently studied. A neural network was used by Kopitkov and Indelman to learn a viewpoint-dependent measurement model of CNN classifier output features [15]. That study focused on the learning procedure of both viewpoint-dependent and class-dependent semantic model, and how to incorporate it as a factor in a factor graph framework. Feldman and Indelman [8] learned a viewpoint-dependent model of semantic measurements, and represent the model as a Gaussian process. Tchuiev and Indelman [16] developed a method for sequential reasoning about the posterior uncertainty of a semantic model. This method was later applied to autonomous planning scenarios [17]. Feldman and Indelman [18] replaced categorical class variables with latent continuous object representations. This representation provides the potential to represent semantic information in a more meaningful way than only the object's class.

Bowman et al. [6] showed that a joint hybrid belief of trajectory, data association, landmark locations, and landmark class can be maximized through expectation-maximization (EM). In the case of scenarios containing many classes and objects, the number of hypotheses through which the expectation must pass is too large and therefore is still infeasible. Tchuiev et al. [7] used a hybrid belief to infer also data association. Data association hypotheses with small probabilities were pruned. They used a semantic model within the factor graph, improving both classification and localization. However, there is no guarantee that the updated belief after pruning will be close to the original belief. Furthermore, due to the high computational complexity of using such a model, they were limited to scenarios with very few classes and objects.

## III. PROBLEM FORMULATION AND NOTATIONS

Suppose a robot travels in an unknown environment and receives observations from various objects as it moves. Let $N$ denote the number of objects the robot observes and $M$ the number of possible classes. Denote the robot's pose at time-step $k$ by $x_k$, and its trajectory from the beginning to time-step $k$ by $x_{1:k} = \{x_1, \ldots, x_k\}$. Let $x_n^o$ represent the pose of the $n$th object, and $X^o = \{x_1^o, \ldots, x_N^o\}$ the poses of all $N$ objects. The class of the $n$th object is $c_n$, and the concatenation of all objects' classes is $C = \{c_1, \ldots, c_N\}$.

The robot receives observations $z_k = \{z_{k,1}, \ldots, z_{k,N_k}\}$ at the $k$th time step, where $N_k$ is the number of objects observed at time-step $k$. In general, it is unknown which observation $z_{k,j}$ corresponds to which of the objects. Let the data association (DA) random variable $\beta_{k,j} = n$ represent the event that observation $z_{k,j}$ corresponds to the $n$th object, and, $\beta_k = \{\beta_{k,1}, \ldots, \beta_{k,N_k}\}$ is the concatenation of all DAs at time-step $k$. In this paper, we assume that DA is known and will extend to unknown DA in the future. Each observation $z_{k,j}$ consists of a geometric part $z_{k,j}^g$ and a semantic part $z_{k,j}^s$, that are assumed to be independent on each other. The observations model is given by

$$
\mathbb{P}\left(z_{k,j} \mid x_{\beta_{k,j}}^o, x_i, c_{\beta_{k,j}}\right) =
$$
$$
\mathbb{P}\left(z_{k,j}^s \mid x_{\beta_{k,j}}^o, x_k, c_{\beta_{k,j}}\right) \mathbb{P}\left(z_{k,j}^g \mid x_{\beta_{k,j}}^o, x_k\right), \quad (1)
$$

where $\mathbb{P}\left(z_{k,j}^s \mid x_{\beta_{k,j}}^o, x_k, c_{\beta_{k,j}}\right)$ is the viewpoint-dependent semantic model, and $\mathbb{P}\left(z_{k,j}^g \mid x_{\beta_{k,j}}^o, x_k\right)$ is the geometric model, and both are assumed to be given. Additionally, the actions $a_{0:k-1}$ and the motion model $\mathbb{P}\left(x_k \mid x_{k-1}, a_{k-1}\right)$ are also assumed to be given.

Define $\mathcal{X}_k = \{X^o, x_{1:k}\}$ as the concatenation of all unknown continuous random variables; we shall refer to it as the state. Define history at time-step $k$ as $H_k = \{z_{1:k}, a_{1:k-1}\}$. The joint hybrid belief over $\mathcal{X}_k$ and $C$ is defined as follows

$$
b\left[\mathcal{X}_k, C\right] \triangleq \mathbb{P}\left(\mathcal{X}_k, C \mid H_k\right). \quad (2)
$$

We present a recursive derivation of the belief. Based on Bayes' theorem, the last observation can be pulled from the history

$$
\mathbb{P}\left(\mathcal{X}_k, C \mid H_k\right) = \eta_k \mathbb{P}\left(z_k \mid \mathcal{X}_k, C, H_k^-\right) \mathbb{P}\left(\mathcal{X}_k, C \mid H_k^-\right), \quad (3)
$$

where $\eta_k \triangleq \mathbb{P}\left(z_k \mid H_k^-\right)^{-1}$ is the normalization factor and $H_k^-$ is the the history without the last observations, $z_k$. The observations are dependent only on the state at the current time, therefore

$$\mathbb{P}\left(z_k \mid \mathcal{X}_k, C, H_k^-\right) = \prod_{j=1}^{N_k} \mathbb{P}\left(z_{k,j} \mid x_k, c_{\beta_{k,j}}, x_{\beta_{k,j}}^o\right) \quad (4)$$

Applying chain rule, the last term in (3) can be written as

$$\mathbb{P}\left(\mathcal{X}_k, C \mid H_k^-\right) = \mathbb{P}\left(x_k \mid a_{k-1}, x_{k-1}\right) \mathbb{P}\left(\mathcal{X}_{k-1}, C \mid H_{k-1}\right). \quad (5)$$

This process can be repeated recursively, resulting in the following formulation

$$\tilde{b}_k\left[\mathcal{X}_k, C\right] = \mathbb{P}_0\left(C\right) \mathbb{P}_0\left(\mathcal{X}_0\right) \prod_{i=1}^{k} \mathbb{P}\left(x_i \mid a_{i-1}, x_{i-1}\right)$$
$$\times \prod_{j=1}^{N_k} \mathbb{P}\left(z_{i,j} \mid x_i, c_{\beta_{i,j}}, x_{\beta_{i,j}}^o\right), \quad (6)$$

where $c_{\beta_{i,j}}$ is the element of $C$ corresponding to index $\beta_{i,j}$. By normalizing (6) we get the belief which is the joint probability of $\mathcal{X}_k$ and $C$

$$b_k\left[\mathcal{X}_k, C\right] = \eta_{1:k}\tilde{b}_k\left[\mathcal{X}_k, C\right], \quad (7)$$

where the normalization factor $\eta_{1:k}$ is given by

$$\eta_{1:k}^{-1} = \sum_{C \in \mathcal{C}} \int_{\Omega_{\mathcal{X}_k}} \tilde{b}_k\left[\mathcal{X}_k, C\right] d\mathcal{X}_k. \quad (8)$$

The belief can be rearranged so that all observations related to the same object are grouped together. Consider $I(n, k)$ as the set of indices $\langle i, j \rangle$ of observations associated with the $n$th object that were observed up to time-step $k$. $I(n, k)$ can be defined mathematically as follows

$$I(n, k) = \{\langle i, j \rangle : \beta_{i,j} = n, 1 \leq i \leq k, 1 \leq j \leq N_i\}. \quad (9)$$

The belief can then be reordered as follows

$$b\left[\mathcal{X}_k, C\right] = \eta_{1:k}\mathbb{P}_0\left(C\right) \mathbb{P}_0\left(X^0\right) \left(\prod_{i=1}^{k} \mathbb{P}\left(x_i \mid a_{i-1}, x_{i-1}\right)\right)$$
$$\times \prod_{n=1}^{N} \prod_{i,j \in I(n,k)} \mathbb{P}\left(z_{i,j} \mid x_n^o, x_i, c_n\right). \quad (10)$$

Using a viewpoint-dependent semantic model, it is evident that classes of different objects are dependent. There are two reasons for this. First, the prior probability of the classes may be dependent, i.e., $\mathbb{P}_0\left(C\right) \neq \prod_{n=1}^{N} \mathbb{P}_0\left(c_n\right)$, for example, near crosswalks we expect to see pedestrians. Second, because of the viewpoint-dependent semantic model, the classes are dependent on the poses, which are in turn dependent on each other, so the classes are dependent.

In general, to infer $\mathcal{X}_k$ and $C$, it is necessary to go through all possible combinations of classes, even when using an efficient algorithm such as expectation maximization. The number of combinations, or hypotheses, is $M^N$. As a result, all algorithms that consider all of the hypotheses run at least at $O\left(M^N\right)$. Because of this, pruning the vast majority of

hypotheses is essential. Using the naive pruning approach, after the pruning and renormalization, it is impossible to determine whether the pruned hypotheses maintain a high probability. In naive pruning, the probabilities of the remaining hypotheses are assumed to sum to one, meaning that the probabilities of the pruned hypotheses are assumed to be zero.

Yet, it is irresponsible to assume that the probabilities of the pruned set are zero. Consider that a robot operates in a human environment, and the robot can do some potentially dangerous things only if it is confident that no human is nearby. By using naive pruning, the robot will be overconfident in the retained hypotheses. In the event that the correct hypothesis was pruned, because there is no indication that the correct hypothesis was pruned, the robot may assume that there is no human nearby and proceed to take the dangerous action. In contrast, if the robot obtains the exact probability that the correct hypothesis is pruned, or at least a lower bound on this probability, the robot will know that it lacks sufficient confidence to take the dangerous action and that the correct hypothesis was pruned.
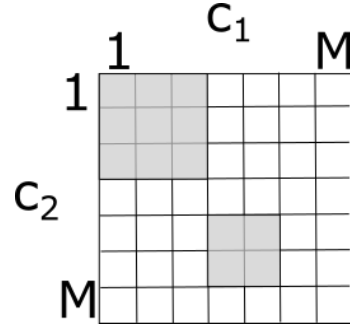


**Fig. 1:** Illustration of the pruned set $\Omega_k^{out}$ in white and the retained set $\Omega_k^{in}$, in gray in a simple case of two objects. A hypothesis consists of a class for each object. It is considered that $|\Omega_k^{in}| \ll |\Omega_k^{out}|$.

## IV. APPROACH

We propose two alternative approaches to evaluating the normalization factor after pruning. We consider two cases. In the first case, the prior probability of the classes is independent, i.e. $\mathbb{P}_0\left(C\right) = \prod_{n=1}^{N} \mathbb{P}_0\left(c_n\right)$. Our study will demonstrate that in this case the original normalization factor can be calculated with the same computational efficiency as the pruning version of the normalization factor. Because we have the original normalization factor, we can query the original probability of each hypothesis. In this case, we know the real probabilities of the remaining hypotheses, and if they are low, we should replace them.

In the second case, the prior probability of the classes is dependent, i.e. $\mathbb{P}_0\left(C\right) \neq \prod_{n=1}^{N} \mathbb{P}_0\left(c_n\right)$. In this case, we propose an upper bound on the probabilities of the pruned hypotheses. Using this bound, we derive a lower bound on the normalization factor. Using an upper bound on the probability of the pruned hypotheses, we can determine whether there is a high probability hypothesis within the pruned set. This bound is also computed with the same computational efficiency.

In both cases, we denote the set of maintained hypotheses by $\Omega_k^{in}$ and the rest of the (pruned) hypotheses by $\Omega_k^{out}$ (see illustration in Fig. 1).

### A. Independent Class Prior

In order to simplify notation, the unnormalized belief (6) can be divided into two parts, one containing the semantic observations, and the other containing the rest.

$$
\tilde{b}_k\left[\mathcal{X}_k, C\right] =
$$
$$
\left(\mathbb{P}_0\left(X^0\right) \prod_{i=1}^{k} \mathbb{P}\left(x_i \mid a_{i-1}, x_{i-1}\right) \prod_{j=1}^{N_i} \mathbb{P}\left(z_{i,j}^g \mid x_{\beta_{i,j}}^o, x_i\right)\right)
$$
$$
\times \left(\mathbb{P}_0\left(C\right) \prod_{i=1}^{k} \prod_{j=1}^{N_i} \mathbb{P}\left(z_{i,j}^s \mid x_{\beta_{i,j}}^o, x_i, c_{\beta_{i,j}}\right)\right). \quad (11)
$$

Define the geometric unnormalized belief as follows

$$
\tilde{b}_k^g\left[\mathcal{X}_k\right] \triangleq \mathbb{P}_0\left(\mathcal{X}_0\right) \prod_{i=1}^{k} \mathbb{P}\left(x_i \mid a_{i-1}, x_{i-1}\right) \prod_{j=1}^{N_i} \mathbb{P}\left(z_{i,j}^g \mid x_{\beta_{i,j}}^o, x_i\right)
$$
$$
(12)
$$

where $H_k^g$ is the history without the semantic observations. The unnormalized belief is equal to the probability of $\mathcal{X}_k \mid H_k^g$ multiplied by the normalization factor $\eta_{1:k}^g$. By substituting (12) into (11), and reordering the equation as in (10), one obtains

$$
\tilde{b}_k\left[\mathcal{X}_k, C\right] = \tilde{b}_k^g\left[\mathcal{X}_k\right] \mathbb{P}_0\left(C\right) \prod_{n=1}^{N} \prod_{i,j \in I(n,k)} \mathbb{P}\left(z_{i,j}^s \mid x_n^o, x_i, c_n\right).
$$
$$
(13)
$$

Further simplifying the belief, we define $\psi_k$ as follows

$$
\psi_k(n, c, \mathcal{X}_k) \triangleq \prod_{i,j \in I(n,k)} \mathbb{P}\left(z_{i,j}^s \mid x_i, x_n^o, c_n = c\right). \quad (14)
$$

It is the likelihood of the object's class $c_n = c$ and the state, given the semantic observations associated with the $n$th object. Substituting $\psi_k$ into (13) results in

$$
\tilde{b}_k\left[\mathcal{X}_k, C\right] = \tilde{b}_k^g\left[\mathcal{X}_k\right] \mathbb{P}_0\left(C\right) \prod_{n=1}^{N} \psi_k(n, c_n, \mathcal{X}_k). \quad (15)
$$

The normalization factor is given in (8), and by substituting (15), we obtain

$$
\eta_{1:k}^{-1} = \sum_{C \in \mathcal{C}} \int_{\Omega_{\mathcal{X}_k}} \tilde{b}_k^g\left[\mathcal{X}_k\right] \mathbb{P}_0\left(C\right) \prod_{n=1}^{N} \psi_k(n, c_n, \mathcal{X}_k) d\mathcal{X}_k. \quad (16)
$$

The above integration, (16), can be represented as expectation over $\mathcal{X}_k \mid H_k^g$ up to the scale of the normalization factor $\eta_{1:k}^g$ of the geometric belief, results in the following

$$
\eta_{1:k}^{-1} = (\eta_{1:k}^g)^{-1} \sum_{C \in \mathcal{C}} \mathbb{P}_0\left(C\right) \mathbb{E}_{\mathcal{X}_k \mid H_k^g}\left[\prod_{n=1}^{N} \psi_k(n, c_n, \mathcal{X}_k)\right].
$$
$$
(17)
$$

Since the normalized belief is required only for the probabilities of classes, the state will be marginalized later. Thus, we could apply the same expectation as in (17), and the

normalization factor $\tilde{b}_k^g$ will be eliminated. Until now, the derivation has been general and does not assume independent priors and therefore applies also to the dependent case.

By using the independent property of the prior of $C$, and reorganizing (16), we obtain

$$
\eta_{1:k}^{-1} = \mathbb{E}_{\mathcal{X}_k \mid H_k^g}\left[\sum_{c_1=1}^{M} \cdots \sum_{c_N=1}^{M} \prod_{n=1}^{N} \mathbb{P}_0\left(c_n\right) \psi_k(n, c_n, \mathcal{X}_k)\right].
$$
$$
(18)
$$

Our *key observation* is that it is possible to substitute the sum and the product, rewriting (18) as

$$
\eta_{1:k}^{-1} = \mathbb{E}_{\mathcal{X}_k \mid H_k^g}\left[\prod_{n=1}^{N} \mathbb{P}_0\left(c_n = c\right) \sum_{c=1}^{M} \psi_k(n, c, \mathcal{X}_k)\right], \quad (19)
$$

since each $c_n$ is independent of the other, inside the expectation. Using (19), the computing complexity of the normalization factor is significantly reduced.

Let us compare the running time of the naive approach in (17) and the efficient approach in (19). In both cases, integration of the state and calculation of $\psi_k(n, c, \mathcal{X}_k)$ is required. Since this calculation is common, it is not included in the analysis to follow. Consider that the expectation over $\mathcal{X}_k$ is approximated by $N_s$ samples drawn from $b_k^g[\mathcal{X}_k]$. Using the naive approach, one must sum the unnormalized probabilities for all hypotheses, as shown in (17). Therefore, the naive approach runs in $O(M^N \cdot N_s)$. In contrast, the computational complexity of the efficient approach (19) is $O\left(N \cdot M \cdot N_s\right)$.

To summarize, using the exact normalization factor (17), the *exact* probability can be computed for each of the maintained hypotheses in $\Omega_k^{in}$. Moreover, the exact probabilities of hypotheses in $\Omega_k^{in}$, provide the exact probability of $\mathbb{P}(C \in \Omega_k^{out})$, which can indicate if the correct hypothesis was pruned. In Fig. 1, and illustration of $\Omega_k^{in}$ and $\Omega_k^{out}$ is provided.

### B. Dependent Class Prior

In this case we consider the prior of the classes is dependent. Suppose we have an upper bound $\tilde{U}_k$ that fulfills the following inequality

$$
\sum_{C \in \Omega_k^{out}} \tilde{b}_k\left[C\right] \leq \tilde{U}_k. \quad (20)
$$

The normalization factor can be bounded from below as

$$
\eta_{1:k} = \left(\sum_{C \in \mathcal{C}} \tilde{b}_k\left[C\right]\right)^{-1} \geq \left(\sum_{C \in \Omega_k^{in}} \tilde{b}_k\left[C\right] + \tilde{U}_k\right)^{-1} \triangleq \underline{\eta}_{1:k},
$$
$$
(21)
$$

where $\underline{\eta}_{1:k}$ is the normalization factor induced by the bound. The probabilities of $C$ can be bounded from below

$$
b\left[C\right] = \eta_{1:k} \tilde{b}\left[C\right] \geq \underline{\eta}_{1:k} \tilde{b}\left[C\right], \quad (22)
$$

A lower bound for the class's probability is a significant result. It enables an autonomous robot to take cautious actions, and indicates whether $\mathbb{P}(C \in \Omega_k^{out})$ is high.
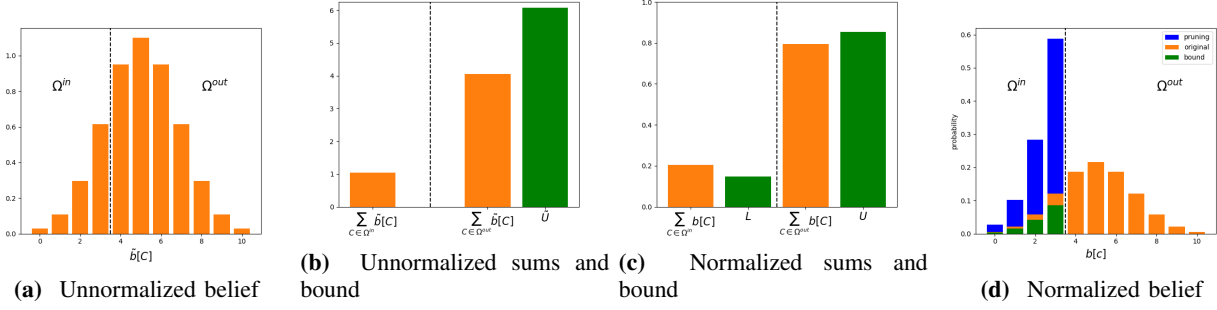
**(a)** Unnormalized belief

**(b)** Unnormalized sums and bound

**(c)** Normalized sums and bound

**(d)** Normalized belief

Fig. 2: The bound process is illustrated in the following figures. (2a) Unnormalized belief separated to $\Omega^{in}$ in the left (before the dashed line) and $\Omega^{out}$ in the right. In (2b), $\sum_{C \in \Omega^{in}} \tilde{b}[C]$ and $\tilde{U}$ are computed. $\sum_{C \in \Omega^{out}} \tilde{b}[C]$ is not computed. In (2c), the sums bound are normalized using $\eta$. A comparison is shown (2d) between the original normalized probabilities, the bound and the naive pruning approach. Using naive pruning can lead to absurd overconfident probabilities. Furthermore, it cannot indicate that a better hypothesis is in the pruned set $\Omega^{out}$. In contrast, if there is a high probability that the correct hypothesis was pruned, then using the upper bound $U$ will provide an indication of this and will allow for a more conservative and realistic normalization of the retained hypotheses. Suppose a hypothesis in $\Omega^{in}$ has probability $\eta \tilde{b}[C]$ higher than $U$, then we are guaranteed that this hypothesis has higher probability than any hypothesis in $\Omega^{out}$.

Furthermore, we can also bound the probability $\mathbb{P}(C \in \Omega_k^{out} \mid H_k)$ from above.

$$\mathbb{P}\left(C \in \Omega_k^{out} \mid H_k\right) = \sum_{C \in \Omega_k^{out}} b[C]$$
$$= 1 - \sum_{C \in \Omega_k^{in}} b[C] \le 1 - \sum_{C \in \Omega_k^{in}} \eta_{1:k} \tilde{b}[C]. \quad (23)$$

This bound is also useful, since it indicates if it is likely that the correct hypothesis was pruned. Fig. 2 illustrates the derivation of the bound clearly.

Next, we will derive the bound and present a method for updating it efficiently. We begin with the simpler version, where the state is assumed to be known, and then we generalize it to the case where it is unknown.

*1) Derivation of $\tilde{U}_k$:* We will begin by considering that the $\mathcal{X}_k$ is known and then extend it to an unknown $\mathcal{X}_k$. The unnormalized belief of $C$ can be derived from degenerating (6)

$$\tilde{b}_k[C] = \mathbb{P}_0(C) \prod_{i=1}^{k} \prod_{j=1}^{N_k} \mathbb{P}\left(z_{i,j}^s \mid x_i, c_{\beta_{i,j}}, x_{\beta_{i,j}}^o\right). \quad (24)$$

We obtain the following belief by reorganizing (24) as was done in (10), and plugging in $\psi_k(n, c, \mathcal{X}_k)$

$$\tilde{b}_k[C] = \mathbb{P}_0(C) \prod_{n=1}^{N} \psi_k(n, c_n, \mathcal{X}_k). \quad (25)$$

According to (20), we need to bound (26) from above

$$\sum_{C \in \Omega_k^{out}} \tilde{b}_k[C] = \sum_{C \in \Omega_k^{out}} \mathbb{P}_0(C) \prod_{n=1}^{N} \psi_k(n, c, \mathcal{X}_k). \quad (26)$$

Using the Cauchy-Schwarz inequality, (26) can be bounded from above as follows

$$\sum_{C \in \Omega_k^{out}} \tilde{b}_k[C] \le \left(\sum_{C \in \Omega_k^{out}} \mathbb{P}_0^2(C)\right)^{1/2} \left(\sum_{C \in \Omega_k^{out}} \prod_{n=1}^{N} \psi_k^2(n, c, \mathcal{X}_k)\right)^{1/2}. \quad (27)$$

Furthermore, we can use the Hölder's inequality, which is an extension of the Cauchy–Schwarz inequality. The Hölder's inequality states that for any two vectors $u, v$ in some inner product space, the following inequality holds

$$|\langle u, v \rangle| \le \|u\|_{q_1} \cdot \|v\|_{q_2}, \quad (28)$$

for any $q_1, q_2 \ge 1$ satisfy $\dfrac{1}{q_1} + \dfrac{1}{q_2} = 1$, where $\|\cdot\|_q$ is the $L(q)$ norm. Therefore,

$$\tilde{U}_k(\mathcal{X}_k) = \left(\sum_{C \in \Omega_k^{out}} \mathbb{P}_0^{q_1}(C)\right)^{1/q_1} \left(\sum_{C \in \Omega_k^{out}} \prod_{n=1}^{N} \psi_k^{q_2}(n, c, \mathcal{X}_k)\right)^{1/q_2}. \quad (29)$$

*2) Efficient update of $\tilde{U}_k$:* Two events can cause an update: a new observation is received, and a change in the composition of the $\Omega_k^{in}, \Omega_k^{out}$. In the former case, the bound should be updated incrementally like the belief. In the latter case, we avoid starting from scratch and re-use previous calculations.

Suppose that the number of elements in $\Omega_k^{in}$ is limited to $N_{in} = |\Omega_k^{in}| \ll |\Omega_k^{out}|$. A sum over $\Omega_k^{out}$ can be replaced with sum over $\mathcal{C}$ minus sum over $\Omega_k^{in}$. Using this attribute, the bound can be rewritten as follows

$$\tilde{U}_k(\mathcal{X}_k) = \left(\sum_{C \in \mathcal{C}} \mathbb{P}_0^{q_1}(C) - \sum_{C \in \Omega_k^{in}} \mathbb{P}_0^{q_1}(C)\right)^{1/q_1}$$
$$\times \left(\sum_{C \in \mathcal{C}} \prod_{n=1}^{N} \psi_k^{q_2}(n, c, \mathcal{X}_k) - \sum_{C \in \Omega_k^{in}} \prod_{n=1}^{N} \psi_k^{q_2}(n, c, \mathcal{X}_k)\right)^{1/q_2}. \quad (30)$$

In order to simplify $\tilde{U}_k$, the sums inside (30) can be defined as following

$$S_{\mathcal{C}}^0 \triangleq \sum_{C \in \mathcal{C}} \mathbb{P}_0^{q_1}(C) \quad (31)$$

$$S_{\Omega_k^{in}}^0 \triangleq \sum_{C \in \Omega_k^{in}} \mathbb{P}_0^{q_1}(C) \quad (32)$$

$$S_{\Omega_k^{in}}^{\psi_k}(\mathcal{X}_k) \triangleq \sum_{C \in \Omega_k^{in}} \prod_{n=1}^{N} \psi_k^{q_2}(n, c_n, \mathcal{X}_k) \quad (33)$$

$$S_{\mathcal{C}}^{\psi_k}(\mathcal{X}_k) \triangleq \sum_{C \in \mathcal{C}} \prod_{n=1}^{N} \psi_k^{q_2}(n, c, \mathcal{X}_k). \quad (34)$$

In (33), the class of the $n$th object is the $n$th element of the vector $C$. In light of definitions (31-34), the bound (30) can

be rewritten as follows

$$\tilde{U}_k(\mathcal{X}_k) = \left(S_{\mathcal{C}}^0 - S_{\Omega_k^{in}}^0\right)^{1/q_1} \left(S_{\mathcal{C}}^{\psi_k}(\mathcal{X}_k) - S_{\Omega_k^{in}}^{\psi_k}(\mathcal{X}_k)\right)^{1/q_2}. \tag{35}$$

*Update due to new observation:* Observing new observation $z_k$ will change $S_{\mathcal{C}}^{\psi_k}$, $S_{\Omega_k^{in}}^{\psi_k}$ and $\psi_k$. As before, we exclude the time complexity of computing $\psi_k$. The dependency on $\mathcal{X}_k$ is omitted here for simplicity of notation, keeping in mind that $\psi_k$ is dependents on $\mathcal{X}_k$. The computation time for the update of $\psi_k$ is $O(M \cdot N_k)$. For the update of $S_{\Omega_k^{in}}^{\psi_k}$, define

$$\varphi_k(C) \triangleq \prod_{n=1}^{N} \psi_k(n, c = c_n, \mathcal{X}_k) \quad C \in \Omega_k^{in}. \tag{36}$$

It is clear that

$$S_{\Omega_k^{in}}^{\psi_k} = \sum_{C \in \Omega_k^{in}} \varphi_k(C). \tag{37}$$

$\varphi_k(C)$ updates as follows

$$\varphi_k(C) = \varphi_{k-1}(C) \prod_{j=1}^{N_k} \mathbb{P}^{q_2}\left(z_{k,j}^s \mid x_{\beta_{k,j}}^o, x_k, c = c_{\beta_{k,j}}\right). \tag{38}$$

The computational time for updating $\varphi_k(C)$ is $O(N_k)$, it should be done for all $C \in \Omega_k^{in}$, therefore, the computation time for updating $S_{\Omega_k^{in}}^{\psi_k}$ is $O(N_k N_{in})$.

The explicit formulation of $S_{\mathcal{C}}^{\psi_k}$ is given by

$$S_{\mathcal{C}}^{\psi_k} = \sum_{C \in \mathcal{C}} \prod_{n=1}^{N} \psi_k(n, c, \mathcal{X}_k)$$
$$= \sum_{c_1=1}^{M} \cdots \sum_{c_N=1}^{M} \prod_{n=1}^{N} \psi_k(n, c = c_n, \mathcal{X}_k). \tag{39}$$

The sum and the product can be substituted for each other in this case

$$S_{\mathcal{C}}^{\psi_k} = \prod_{n=1}^{N} \sum_{c_n=1}^{M} \psi_k(n, c = c_n, \mathcal{X}_k). \tag{40}$$

Define

$$s_n^{\psi_k} \triangleq \sum_{c=1}^{M} \psi_k(n, c, \mathcal{X}_k). \tag{41}$$

(41) can be inserted to $S_{\mathcal{C}}^{\psi_k}$

$$S_{\mathcal{C}}^{\psi_k} = \prod_{n=1}^{N} s_n^{\psi_k}. \tag{42}$$

The update affects only $s_n^{\psi_k}$ for $n \in \beta_k$. The computation time for updating $s_n^{\psi_k}$ is $O(M)$ and for $S_{\mathcal{C}}^{\psi_k}$ is $O(M \cdot N_k)$. The final computation time for an observation-based update is $O(N_k \max(M, N_{in}))$.

*Update due to change in* $\Omega_k^{in}, \Omega_k^{out}$: The second type of update is caused by a change in the composition of $\Omega_k^{in}, \Omega_k^{out}$, and affects only $S_{\Omega_k^{in}}^0, S_{\Omega_k^{in}}^{\psi_k}$. Consider adding a new hypothesis, $C^a$, and removing existing hypothesis $C^r$ from $\Omega_k^{in}$. The update will be performed as follows

$$S_{\Omega_k^{in}}^0 \leftarrow S_{\Omega_k^{in}}^0 + \mathbb{P}_0^{q_1}(C^a) \tag{43}$$

$$\varphi(C^a) = \prod_{n=1}^{N} \psi_k(n, c = c_n^a) \tag{44}$$

$$S_{\Omega_k^{in}}^{\psi_k} \leftarrow S_{\Omega_k^{in}}^{\psi_k} + \varphi(C^a). \tag{45}$$

The update involves computing $\varphi(C^a)$, which is $O(N)$. Now the removal of $C^r$ is done by

$$S_{\Omega_k^{in}}^0 \leftarrow S_{\Omega_k^{in}}^0 - \mathbb{P}_0^{q_1}(C^r) \tag{46}$$

$$S_{\Omega_k^{in}}^{\psi_k} \leftarrow S_{\Omega_k^{in}}^{\psi_k} - \varphi(C^r), \tag{47}$$

which is done in $O(1)$.

### C. Belief over unknown classes and state

The unnormalized belief over the classes and the state can be taken from (15). According to (20), $\tilde{U}_k$ should satisfy

$$\tilde{U}_k \geq \sum_{C \in \Omega_k^{out}} \int_{\Omega_{\mathcal{X}_k}} \tilde{b}_k^g[\mathcal{X}_k] \, p_0(C) \prod_{n=1}^{N} \psi_k(n, c_n, \mathcal{X}_k) d\mathcal{X}_k. \tag{48}$$

Integration can be replaced by expectation over $\mathcal{X}_k \mid H_k^g$, as done in (18),

$$\tilde{U}_k \geq \mathbb{E}_{\mathcal{X}_k \mid H_k^g}\left[\sum_{C \in \Omega_k^{out}} p_0(C) \prod_{i=1}^{k} \prod_{j=1}^{n_i} \psi_k(n, c_n, \mathcal{X}_k)\right]. \tag{49}$$

The term inside the expectation of (49) is the unnormalized belief of $C \mid \mathcal{X}_k$ given in (25). Therefore, it can also be bounded using the bound in (35), thus

$$\tilde{U}_k = \mathbb{E}_{\mathcal{X}_k \mid H_k^g}\left[\left(S_{\mathcal{C}}^0 - S_{\Omega_k^{in}}^0\right)^{1/q_1}\left(S_{\mathcal{C}}^{\psi_k}(\mathcal{X}_k) - S_{\Omega_k^{in}}^{\psi_k}(\mathcal{X}_k)\right)^{1/q_2}\right]. \tag{50}$$

$\left(S_{\mathcal{C}}^0 - S_{\Omega_k^{in}}^0\right)^{1/q_1}$ is not dependent on $\mathcal{X}_k$ so it can be taken out of the expectation. The time complexity of computing the bound is the same as before multiplied by the number of samples $N_s$.

## V. EXPERIMENTS

In this section, we evaluate our methods using synthetic simulations. It is important to note that we do not propose a new viewpoint-dependent model, but only suggest that given such a model is utilized, our method can reduce computational complexity substantially

A 2D environment is assumed with a robot that moves around, observing the objects that it encounters along the way. Each object has orientation and class type. Objects are randomly positioned throughout the environment, and their classes are randomly selected. Initially, the robot does not know where the objects are placed or to which class they

belong. The geometric observations are bearing measurements and the semantic observations are synthetic classifier output. Based on the geometric measurements and the motion model, the geometric belief is evaluated. The geometric belief is assumed to be Gaussian. To calculate the full SLAM graph, we used the GTSAM library with Python wrapper, resulting in a mean and a covariance of the geometric belief. Throughout all the simulations, 100 samples are drawn from the geometric belief.

For the robot's prior pose, process noise, and geometric observation noise, the covariances were $\Sigma_0 = \text{diag}(0.01, 0.01, 0.001)$, $\Sigma_p = \text{diag}(0.3, 0.3, 0.03)$, and $\Sigma_m = \text{diag}(0.03)$, respectively. The dependent prior probability of $C$ is randomly selected using the following procedure. A matrix $A$ of size $\underbrace{M \times M \cdots \times M}_{N}$ was allocated. The elements of $A$ was sampled from the uniform distribution $U[0.001, 1]$. Then $A$ was divided by the sum of its elements. Thus, the prior probability of $C$ is given by $\mathbb{P}_0(C) = A[C]$. In the case of an independent prior, this procedure was used for each marginal.

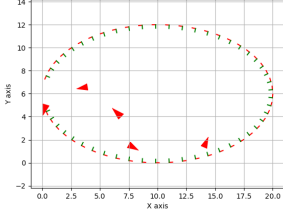In Fig. 3 we can see the scheme of the scene.



**Fig. 3:** Visualization of the scene. The triangles in red represent objects that are oriented. The angles and classes were chosen at random. The robot moves counter-clockwise, starting from the lower left side.

The distribution of the semantic observations is logit-normal distribution, thus,

$$z_{k,j} \sim LN\left(\mu\left(x^o_{\beta_{k,j}}, x_k, c\right), \Sigma\right), \quad (51)$$

where $\Sigma = I_M \sigma_s^2$ and $\sigma^s = 0.015$. The same distribution is used in equation [17] as well. It is possible to interpret this distribution as the output of a classifier since the samples are vectors of positive elements that sum to 1. For simplicity of notation, consider that $z_{k,j}$ is associated with the $n$th object, thus $\beta_{k,j} = n$. The semantic measurement model is defined as follows

$$\mu\left(x^o_n, x_k, c\right) = \begin{cases} e_c \cdot h\left(x^o_n, x_k\right) & c \in [1, M-1] \\ 0_{M-1} & c = M \end{cases}, \quad (52)$$

$$h\left(x^o_n, x_k\right) = (1 - \cos(\theta)) \min\left\{\frac{1}{dist(x_k, x^o_n)}, \frac{1}{2}\right\}, \quad (53)$$

where $\mathbf{e}_c \in \mathbb{R}^{M-1}$ is a vector of zeros with a single element at the $c$th position equal to one, and $dist(\cdot, \cdot)$ is the euclidean distance. $\theta$ is the relative angle, calculated from the relative pose $x^{rel} = x_k \ominus x^o_n$. The mean $\mu\left(x^o_n, x_k, c\right)$ is characterized by the following features: under the the hypothesis $c = M$, the mean of $z_{k,j}$ is $\frac{1}{M}$. For the other hypotheses, as

the camera gets closer to the object and becomes oriented towards the front of the object, the probability that the right element of $z_{k,n}$ will respond to object's class is increased.

We compare between the original probability without pruning, the naive pruning approach, the bound approach, and pruning using the original normalization in the independent prior case.

The sets $\Omega_k^{in}, \Omega_k^{out}$ are predefined and are constant. The size of $\Omega_k^{in}$ is $N_{in} = 8$. A true hypothesis $C^{(t)} = (c_1^{(t)}, \ldots, c_N^{(t)})$ is predefined as well. The semantic observations are sampled from the probability $\mathbb{P}\left(z_{k,n} \mid x^o_n, x_k, c_n^{(t)}\right)$ described in (51), (52), and (53).

The results are divided into two cases: when $C^{(t)} \in \Omega_k^{in}$ and where $C^{(t)} \in \Omega_k^{out}$. In order to determine whether the system is confident in its best hypothesis, we look at the maximum probability in $\Omega_k^{in}$, $\max_{C \in \Omega_k^{in}} \{b[C]\}$, over time. The most likely hypothesis is the same across all methods because the difference between them is only the normalization factors, but different probability is provided by each approach. The same sets $\Omega_k^{in}, \Omega_k^{out}$ used for all methods.

Both Fig. 5a and Fig. 5b illustrate the scenario of classes with dependent priors. The naive exact belief, naive pruning, and bound are compared. The number of objects is 5 and number of classes is 3, the total number of hypotheses is $3^5 = 243$. Even when considering such a small number of classes and objects, the number of hypotheses is not so small.

Fig. 5a shows the case for $C^{(t)} \in \Omega_k^{in}$. The naive pruning approach is very confident in its hypothesis from the beginning to the end. The naive exact approach, which provides the exact probability of this hypothesis, is less confident. The bound, which bounds the exact probability from below (22), following the exact approach throughout the simulation. Although the bound is a lower bound on the exact belief, it remains close to the exact belief even when the latter reaches a high probability.

Fig. 5b, shows the case for $C^{(t)} \in \Omega_k^{out}$. Throughout the simulation, the exact probability is close to zero and therefore also the bound. By contrast, the naive pruning approach is too confident, assuming that the correct hypothesis is in $\Omega_k^{in}$.

The independent prior case is considered in figures Fig. 4a and Fig. 4b. For the independent prior case, we used $N = 5$ and $M = 100$. In this scenario the number of hypotheses is $10^{10}$. In this case, we compare the efficient exact belief with naive pruning and the bound approach. Despite the fact that the bound is not required in this case, we wish to evaluate how it will work when many classes are involved.

In this case, the results appear to be the same as in the dependent case. In Fig. 4a, the naive approach is overconfident. In the exact belief we see that only after several time-steps the belief gain confidence, and the bound is follows the exact belief. According to Fig. 4b, the naive approach is absurdly overconfident, while the exact belief and the bound approach indicate that the probability is zero.

In figures Fig. 6a and Fig. 6b we compare the runtime of naive pruning, naive exact, the efficient exact and the bound approach. For each simulation, one hundred trails were performed and the mean value of the runtime was taken.
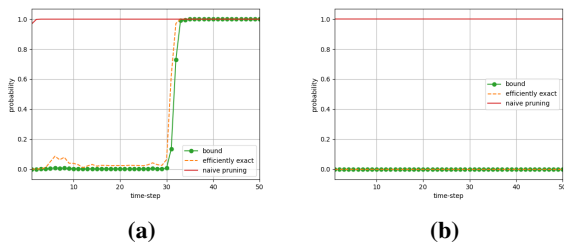
**Fig. 4:** The maximal probability of the maintained hypotheses. Considering an *independent* prior, where **(a)** $C^{(t)} \in \Omega_k^{in}$, and **(b)** $C^{(t)} \in \Omega_k^{out}$.
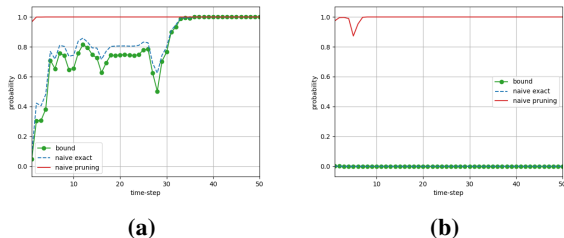


**Fig. 5:** The maximal probability of the maintained hypotheses. Considering an *dependent* prior, where **(a)** $C^{(t)} \in \Omega_k^{in}$, and **(b)** $C^{(t)} \in \Omega_k^{out}$.

In Fig. 6a, show the runtime versus the number of objects $N$. The number of classes is $M = 2$. It appears that the efficient exact and the bound approach have the same runtime as the naive pruning approach. In Fig. 6b, the runtime is plotted against the number of classes $M$. The number of objects is $N = 3$. In comparison to a naive pruning approach, there is a slight increase in runtime for the efficient exact and bound approaches. Given the significant improvement in accuracy and reliability, we consider the efficient exact and bound approaches to be worthwhile.

## VI. CONCLUSIONS

In the context of semantic SLAM, we explore the pruning procedure often used when the size of the alphabet of the hybrid belief is exponentially large, in the case where the semantic observation model is dependent on the relative position of the object. In this scenario, the number of hypotheses is exponential in the number of objects, and therefore, not feasible in a real-time framework. Due to this reason, the pruning procedure is used, but following
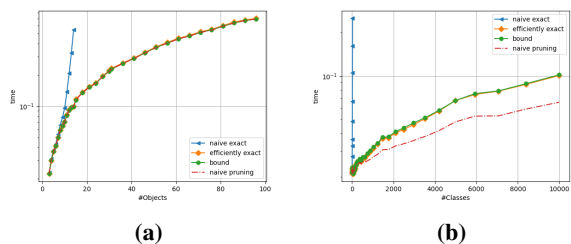


**Fig. 6:** **(a)** Running time versus number of objects $N$. The number of classes is $M = 2$. The computational complexity of both the bound and the efficient method for the independent case original are the same as pruning. **(b)** Running time versus number of classes $M$. The number of object is $N = 3$. Here, the naive pruning running time is slightly better.

pruning and renormalization, the resulting belief is over-confident, disregards the possibility that the true hypothesis is pruned. When the prior over classification variables of different objects is independent, we demonstrated that the normalization factor can be calculated efficiently, allowing us to query the exact probability of each hypothesis individually. In the dependent case, using our method, we bound the probability that the true hypothesis was pruned, and, as a result, we obtain a more accurate and conservative estimate of the original belief, with the guarantee that the original belief over classification variables without pruning is equal or higher than our belief after pruning. Future research might explore the minimization of the upper bound $U_k$ with respect to $\Omega_k^{in}$ as a method to decide which hypotheses to maintain and which to prune.

## REFERENCES

[1] R. Smith, M. Self, and P. Cheeseman, "Estimating uncertain spatial relationships in Robotics," in *Autonomous Robot Vehicles*, I. Cox and G. Wilfong, Eds. Springer-Verlag, 1990, pp. 167–193.

[2] J. Leonard and H. Durrant-Whyte, "Mobile robot localization by tracking geometric beacons," *IEEE Trans. Robot. Automat.*, vol. 7, no. 3, pp. 376–382, 1991.

[3] H. Durrant-Whyte, S. Majunder, S. Thrun, M. de Battista, and S. Scheding, "A Bayesian algorithm for simultaneous localization and map building," in *Proceedings of the 10th International Symposium of Robotics Research*, 2001.

[4] S. Thrun, *Simultaneous Localization and Mapping*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 13–41.

[5] K. Doherty, D. Fourie, and J. Leonard, "Multimodal semantic slam with probabilistic data association," in *2019 international conference on robotics and automation (ICRA)*. IEEE, 2019, pp. 2419–2425.

[6] S. Bowman, N. Atanasov, K. Daniilidis, and G. Pappas, "Probabilistic data association for semantic slam," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 1722–1729.

[7] V. Tchuiev, Y. Feldman, and V. Indelman, "Data association aware semantic mapping and localization via a viewpoint-dependent classifier model," in *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2019.

[8] Y. Feldman and V. Indelman, "Bayesian viewpoint-dependent robust classification under model and localization uncertainty," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2018.

[9] S. Pillai and J. Leonard, "Monocular slam supported object recognition," in *Robotics: Science and Systems (RSS)*, 2015.

[10] S. Yang and S. Scherer, "Cubeslam: Monocular 3-d object slam," *IEEE Transactions on Robotics*, vol. 35, no. 4, pp. 925–938, 2019.

[11] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison, "Slam++: Simultaneous localisation and mapping at the level of objects," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2013, pp. 1352–1359.

[12] S. Omidshafiei, B. T. Lopez, J. P. How, and J. Vian, "Hierarchical bayesian noise inference for robust real-time probabilistic object classification," *arXiv preprint arXiv:1605.01042*, 2016.

[13] L. Nicholson, M. Milford, and N. Sünderhauf, "Quadricslam: Dual quadrics from object detections as landmarks in object-oriented slam," *IEEE Robotics and Automation Letters (RA-L)*, vol. 4, no. 1, pp. 1–8, 2018.

[14] K. J. Doherty, D. P. Baxter, E. Schneeweiss, and J. J. Leonard, "Probabilistic data association via mixture models for robust semantic slam," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 1098–1104.

[15] D. Kopitkov and V. Indelman, "Robot localization through information recovered from cnn classificators," in *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, October 2018.

[16] V. Tchuiev and V. Indelman, "Inference over distribution of posterior class probabilities for reliable bayesian classification and object-level perception," *IEEE Robotics and Automation Letters (RA-L)*, vol. 3, no. 4, pp. 4329–4336, 2018.

[17] ——, "Epistemic uncertainty aware semantic localization and mapping for inference and belief space planning," *arXiv preprint arXiv:2105.12359*, 2021.

[18] Y. Feldman and V. Indelman, "Towards self-supervised semantic representation with a viewpoint-dependent observation model," in *Workshop on Self-Supervised Robot Learning, in conjunction with Robotics: Science and Systems (RSS)*, July 2020.