

Active Multi-View Object Recognition and Online Feature Selection

Christian Potthast, Andreas Breitenmoser, Fei Sha and Gaurav S. Sukhatme

Abstract Adaptive action selection is crucial for object recognition on robots with limited operating capabilities. It offers much desired flexibility in trading off between the costs of acquiring information and making robust and reliable inference under uncertainty. In this paper, we describe an information-theoretic framework that combines and unifies two common techniques: *view planning* for resolving ambiguities and occlusions and *online feature selection* for reducing computational costs. Concretely, our algorithm adaptively chooses two strategies: utilize simple-to-compute features that are the most informative for the recognition task or move to new viewpoints that optimally reduce the expected uncertainties on the identity of the object. Extensive empirical studies have validated the effectiveness of the proposed framework. On a large RGB-D dataset, dynamic feature selection alone reduces the computation time at runtime by five folds, and when combining it with viewpoint selection, we significantly increase the recognition accuracy on average by 8 – 15% absolute, compared to systems that do not use these two strategies. Lastly, we have also demonstrated successfully the effectiveness of the framework on a quadcopter platform with limited operating time.

1 Introduction

Recognizing objects from visual percepts is a challenging task for mobile robots. Conventional approaches for object recognition are tour de force, relying almost exclusively on complex statistical models for classification and heavy engineering of computationally-intensive features [4, 14]. The processing capability of the recognition system, however, is seldom questioned and assumed to be unconstrained. In stark contrast, mobile robots are often underpowered and need to be frugal on resources for computing and reasoning—being able to move is their first priority!

Active object recognition[2, 8] is an appealing paradigm to overcome the challenge. Here, a mobile robot configures and positions its sensor to inspect an object and recognize it as a previously learnt one. Active object recognition enables exploratory actions to account for unforeseen objects or scenes, and offers greater

Christian Potthast · Andreas Breitenmoser · Fei Sha · Gaurav S. Sukhatme
University of Southern California (USC), Department of Computer Science,
Los Angeles, California, 90089
e-mail: potthast@usc.edu (Christian Potthast), andreas.breitenmoser@usc.edu (Andreas Breitenmoser), feisha@usc.edu (Fei Sha), gaurav@usc.edu (Gaurav S. Sukhatme)

flexibility in actively seeking the most useful information for inference under uncertainty.

As such, adaptive action selection can potentially harvest the statistical modeling power of conventional approaches even under the severe (computing) constraints on the mobile robots. The general idea is for the sensors to be controlled such that only data with high information content is collected, *adapted dynamically* to the robot’s knowledge (i.e., the belief state about the identity of the object) as a result of running inference on previously acquired information.

Concretely, two common actions have been explored in previous work: to compute a subset of yet informative features and to control the number of different viewpoints. However, those works do not exploit the full potential of combining the two strategies in concert.

The main idea of our work hinges on the following intuition: *how to reduce overall costs for recognition by relying on very simple features extracted at multiple views in lieu of complex features extracted at a single view?* Specifically, in this paper, we describe an information-theoretic framework that unifies these two strategies: view planning for moving robots to the optimal viewpoints and dynamic feature selection for extracting simple yet informative features. However, unlike standard feature selection, the two strategies are symbiotic in our algorithm. The selected features at any given time depend on the past trajectory of the viewpoints while simultaneously affecting the future viewpoint(s).

We conduct extensive evaluations on a large RGB-D camera dataset of single object scenes. Our results show that dynamic feature selection reduces the computation time by five folds at runtime. Furthermore, when we combine it with viewpoint planning, we increase significantly the recognition accuracy, often as high as 8 – 15%. We further demonstrate the proposed approach through experiments with a quadcopter robot, equipped with a RGB-D camera that collects multiple views by flying around an object. Our results show pose estimates within $\pm 45^\circ$ of the ground truth and object recognition accuracies of over 90%, despite using relatively simple features in rather challenging object recognition experiments.

The remainder of the paper is organized as follows. Sec. 2 discusses related work. We formalize the main idea of active multi-view object recognition in Sec. 3, followed by a detailed description of how to plan optimal views and select features online in Sec. 4. We report empirical evaluations in Sec. 5 and conclude in Sec. 6.

2 Related Work

Works on active perception go back as far as the seminal paper [3]. “Active” can either refer to adjustments of the sensor parameters themselves, e.g., adaptively changing the focal length of a camera, as in [3, 8], or mean that the entire sensor moves together with the underlying platform. The latter is particularly common in robotics, where robots plan new observation locations and viewing directions actively [13, 11, 2]. This is known as robotic view planning [18] and has various applications, such as inspection of underwater structures, object search in indoor environments or detecting objects on a table top [2]. The common goal of these and

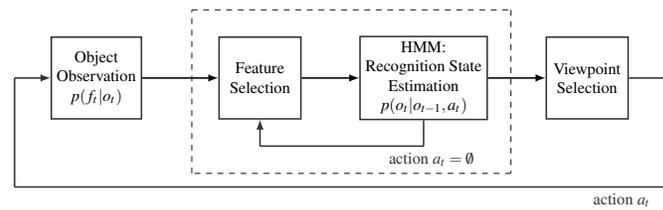


Fig. 1 Overview of the active multi-view object recognition framework. After each observation, the HMM updates the recognition state, upon which a new feature (inner loop) or a new viewpoint (outer loop) is selected. The feature type that is selected first at any new location is pre-defined.

similar approaches is to recognize objects and estimate their pose [8, 13, 12]. To deal with uncertainty in the robot’s state space, [2, 20] formulated the problem as a POMDP, however, their decision making comes with a higher computational cost. [6, 16, 7] propose multi-view object recognition frameworks that use cameras and image databases; alternatively, objects can be recognized by extracting features from laser scanners and RGB-D depth cameras [2, 9, 15]. [15] applies the online boosting approach from [10] to point clouds in the robotic domain for feature selection. Feature selection can support the object recognition task by reducing the computation time and confusions through suboptimal feature sets [21]. A common alternative to boosting is the selection of features by mutual information [17, 21]. [1] evaluates the feature selection methods from [17] when applied to the RGB-D dataset of [12]. Although many approaches for active view planning as well as online feature selection exist, to the best of our knowledge, the advantages of both approaches have never been leveraged in combination. [22] proposes a unifying view on active recognition and feature selection, however, focuses on feature interdependencies and the selection of features given a concrete context, e.g., a specific target object to detect. Their method, as well as the other feature selection methods for object recognition above, do not explicitly consider viewpoint selection, whereas our object recognition framework applies both feature and viewpoint selection in concert.

3 Problem Formulation

Given a static scene and an object in question, we want to recognize the object as well as estimate its pose relative to the coordinate system of the observing sensor. We are interested in choosing a sequence of actions for the sensor which maximizes the recognition accuracy and leads to a minimum operation time and thus minimal energy consumption. In each time step t , we are provided with a set of possible actions a_t and have the choice of 1) stopping the recognition process, 2) staying at the current location and computing a next feature f_t with different feature type (i.e., applying a null action $a_t = \emptyset$), or 3) moving the sensor to a new observation location by taking the action $a_t \neq \emptyset$ and computing a first feature f_t from the new viewpoint. We assume that the type of the first feature evaluated at a new location is pre-defined by the user. Further, we restrict the motion of the sensor to a path that connects fixed discretized viewing positions. This reduces complexity in computation compared to a continuous representation of viewpoints. After an action is taken, we obtain a

new sensor measurement. The sensor measurements are parametrized by features from a feature vector $\mathbf{f} = [f^1, \dots, f^N]$ which includes N different feature types in total. The values of the features depend on the identity of an object. We incorporate the uncertainty of the observations by the conditional probability density function $p(f_i|o_t)$ with $o_t = \{c_t, \theta_t\}$. In order to infer about the current object class c_t and orientation angle θ_t of the observed object from a set of measurements, we use a *Bayesian framework*. The parameters are assumed to be obtained from a physical object model and to be trained offline. The general system diagram of our active recognition framework is presented in Fig. 1. It consists of an inner and an outer loop and the four modules: *object observation*, *feature selection*, *recognition state estimation*, and *viewpoint selection*.

At each observation location, we start off with extracting one first feature. We then have the choice of either extracting another feature (i.e., to locally exploit) or moving to a new location (i.e., to globally explore), depending on our stopping criterion. However, the premise is to extract as much information as possible at the current location since moving to a new location is inherently more costly. Also, we only want to move to a new location if the gain of information (in expectation) significantly reduces the uncertainty of the belief state, or in other words, if the recognition accuracy is expected to improve significantly. Each time we decide to add a new feature, we compute the utility of all remaining features in our arsenal and select the one with the highest utility. We only allow to add one feature at a time, and once a feature is added it cannot be removed from the set anymore. By doing this, integrating a new feature is simply an update step, without the need of recomputing the posterior of the entire feature set. We continue until either we reach our stopping criterion or all N features are observed at a given location, forcing us to move to the next location or to terminate. By this, ideally, at each location we have only extracted useful features and left out features that yield no additional information. Additionally, not extracting a feature saves time, which helps to reduce the total time needed to recognize the object.

Sequential Bayesian Recognition: Given a new feature f_t observed from the object, we want to efficiently integrate the new observation with results from prior observations and update our current belief about the object class and orientation. We formulate the recognition and pose estimation of the object as a Bayesian network in form of a HMM, where the observation sequence is defined as $S = \{o_1, \dots, o_t, a_2, \dots, a_t\}$ as shown in Fig. 2. The hidden state of our HMM is defined as recognition state o_t with $o_{1:T} = \{o_1, \dots, o_t, \dots, o_T\}$, the observations are given as observed features f_t with $f_{1:T} = \{f_1, \dots, f_t, \dots, f_T\}$ and the actions to move between viewpoints are given as a_t with $a_1 = \emptyset$ and $a_{2:T} = \{a_2, \dots, a_t, \dots, a_T\}$. The state of the hidden process satisfies the Markov property, that is, given the value of o_{t-1} , the current state o_t is independent of all the states prior to $t-1$. We can define the joint distribution of the HMM as follows:

$$\begin{aligned} P(o_{1:T}, f_{1:T}, a_{2:T}) &= P(o_{1:T})P(a_{2:T}|o_{1:T})P(f_{1:T}|o_{1:T}, a_{2:T}) \\ &= \left[P(o_1) \prod_{t=2}^T P(o_t|o_{t-1}, a_t) \right] \left[\prod_{t=1}^T P(f_t|o_t) \right] \left[\prod_t P(a_t) \right]. \end{aligned} \quad (1)$$

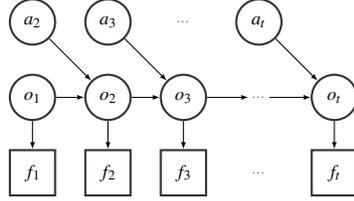


Fig. 2 HMM for sequential Bayesian recognition. Given a sequence of observed features $f_{1:T}$ and executed actions $a_{2:T}$, we reason about the hidden recognition state of an object at time t , $o_t = \{c_t, \theta_t\}$, including its object class $c_t = c^k$ and object pose $\theta_t = \theta^q$. (Null action a_1 is not displayed.)

The *observation model* is given as $P(f_t|o_t)$, representing the likelihood that feature f_t is generated by recognition state o_t . The transition probability is given by $P(o_t|o_{t-1}, a_t)$, which means that the transition to the next recognition state is dependent on the previous recognition state (i.e., the previous object class and orientation) as well as the action that has been applied. The posterior marginals $P(o_t|f_{1:T}, a_{2:T})$ of all hidden state variables given a sequence of observations can be efficiently computed using the *forward-backward algorithm*. The algorithm computes a smoothed estimate given all evidence in two passes, with the first pass going forward in time while the second pass going backwards in time:

$$\begin{aligned} P(o_t|f_{1:T}, a_{2:T}) &= P(o_t|f_{1:t}, f_{t+1:T}, a_{2:t}, a_{t+1:T}) \\ &\propto P(f_{t+1:T}, a_{t+1:T}|o_t) P(o_t|f_{1:t}, a_{2:t}). \end{aligned} \quad (2)$$

In the first pass, we compute the forward probabilities $P(o_t|f_{1:t}, a_{2:t})$, i.e., the probabilities for all $t \in \{1, \dots, T\}$ ending up in a particular state o_t given the first t observations in the sequence. In the second step, we compute the probabilities $P(f_{t+1:T}, a_{t+1:T}|o_t)$ of the remaining observations given any starting point t .

Observation Model: For recognition, we learn a function $g: \mathbb{R}^{|f|} \rightarrow \mathbb{N}$ that maps the extracted object features in \mathbf{f} to a unique object identifier $i \in \{1, \dots, KQ\}$, each pointing to a distinct object model $o^i = \{c^k, \theta^q\}$, with $k \in \{1, \dots, K\}$ and $q \in \{1, \dots, Q\}$. K denotes the total number of object classes c^k and Q defines the number of discretized object poses θ^q per object class. Learning multiple models per object class allows us to train a model for a partial view of an object. This enables inference about an object's pose and reduces the variance in the model. The function is learned in a supervised fashion, with all possible features $\mathbf{f} = [f^1, \dots, f^N]$ and the unique identifier i available as the value pair (\mathbf{f}, i) at model training phase. While we assume that all features of an object are observed at training phase, we do not make this assumption during prediction, where not all components of \mathbf{f} are observed due to feature selection.

To deal with the missing data problem at prediction time, we use a Gaussian generative model of the form $p(f|o^i) = \prod_{m=1}^M \mathcal{N}(f|\mu_{o^i,m}, \sigma_{o^i,m})$. M is the number of independent normal distributions to model a feature f . Each feature f is represented as a feature histogram with M bins, and the mean $\mu_{o^i,m}$ and variance $\sigma_{o^i,m}$ for each bin m are computed over the training data of each object model o^i . We use M inde-

pendent normal distributions instead of a multivariate normal distribution to avoid high-dimensional features, which usually perform rather poorly due to the curse of dimensionality. Defining the problem as an independent Gaussian generative classifier has the additional advantage of handling missing features by marginalizing them out at prediction time, i.e., instead of inferring about the missing features, we simply perform inference on a subset of the full feature vector \mathbf{f} .

The Gaussian representation of the features is compact and fast to train, since it only requires the computation of the Gaussian feature parameters $\mu_{o^i,m}$ and $\sigma_{o^i,m}$. Sequential learning as well as sequential updates of mean and variance are possible, which has the potential of adaptively adding new observations to our trained features. For example, we can compute the observation model for the full feature vector \mathbf{f} as

$$p(\mathbf{f}|o^i) = \prod_{n=1}^N \psi_n \prod_{m=1}^M \mathcal{N}(f^n | \mu_{o^i,m}, \sigma_{o^i,m}), \quad (3)$$

with ψ_n denoting a scaling factor which scales the different feature distributions.

4 Action Selection

In the case of object recognition, the task of action selection is twofold. On the one hand, we would like to select the lowest possible number of viewpoints that make an informed decision about an object possible. On the other hand, however, certain viewing choices or combinations of features may be well explained by more than one hypothesis. This means, the measurement distributions that correspond to the competing hypotheses are similar for a particular viewpoint. With respect to optimization, we want to select the actions that yield the most useful information and to stop acquiring and incorporating new observations once we have reached a desired level of certainty about an object.

Our Bayesian inference framework leads naturally to an iterative information-theoretic approach for deciding on the next action a_{t+1} . Action selection attributes a *utility score* $s_t(a, f)$ to each possible next action a and feature f . Possible next actions are either to stay (i.e., $a = \emptyset$) and compute a next feature f , or to move and compute a new first feature from a different observation location; the observation location is selected from the discrete set of all possible next viewing positions. The next action is eventually evaluated by $f_{t+1} = \operatorname{argmax}_f s_t(a, f)$, for the best feature f while fixing $a = \emptyset$, and by $a_{t+1} = \operatorname{argmax}_a s_t(a, f)$, selecting the best a (without selecting f). Several methods from information theory have been shown to work well in action selection tasks. The more popular methods include *expected loss of entropy* (LE) [16, 6], *mutual information* (MI) [2, 8] and *Kullback-Leibler divergence* (KL) [13], which we adapted in our framework for the action selection task.

Convergence and Stopping Criterion: The sequential update allows to acquire new information and with that reduce the uncertainty of the belief state. Unfortunately, one cannot guarantee that the sequential decision process reduces the uncertainty in every single step. But, on average, by utilizing information-theoretic action selection, the uncertainty will be reduced. The problem with information-theoretic

methods is to find an adequate *stopping criterion*. In the case of feature selection, we need a stopping criterion that defines when to stop adding new features from the current location and rather move to a new location. Similarly, for viewpoint selection we need a stopping criterion to know when we can terminate the object recognition process and rely on the observations and inference we have made so far. The results of an action selection method highly depend on the stopping criterion as well as the parameters controlling that criterion. Wrong choices can lead to sub-optimal performance, with all possible actions a being processed in the worst case.

4.1 Feature Selection

Feature selection can be performed at training time or at prediction time. At training time, the feature space is reduced in advance by selecting relevant features, such that the loss of significant information is minimal. At prediction time, features are selected and computed online. In this work, we train all features and select features only at prediction time. This keeps the complexity low during training and facilitates updates by adding new features without the need to re-train the entire training set.

We want to compute relevant features, i.e., the features should provide additional information and preferably result in a unimodal posterior distribution. This becomes especially important on a robot platform with minimal computation power and short operation time. In the following subsections, we describe two methods for feature selection, one using LE and one using MI.

4.1.1 Feature Selection by Expected Loss of Entropy

Maximum LE indicates a feature that, if extracted next, would lead to the highest reduction in entropy of the posterior distribution on average. In other words, the selection of this feature reduces the uncertainty about our hypothesis on average the most. We compute the average LE because we cannot be certain about which object we are currently observing.

We compute the utility score $s_t(a, f)$ for a null action $a = \emptyset$ and a next possible feature f , based upon the confidence in the current posterior probability $P(o^i|f_{1:t}, a_{2:t})$ obtained from (2), as

$$s_t(a, f) = \sum_{i=1}^{KQ} P(o^i|f_{1:t}, a_{2:t}) E[\Delta H(o^i, f_{1:t}, f, a_{2:t})], \quad (4)$$

with KQ object models in the model set. In this formulation, LE is weighted by the current belief of the object class o^i being the correct hypothesis, and estimated as the expectation $E[\cdot]$ of the difference in entropy between two subsequent posterior distributions,

$$E[\Delta H(o^i, f_{1:t}, f, a_{2:t})] = H(o_t|f_{1:t}, a_{2:t}) - \int p(f|o^i) H(o_t|f_{1:t}, f, a_{2:t}) df. \quad (5)$$

As we do not know which feature to expect, we need to integrate over the entire feature space. The feature distribution $p(f|o^i)$ is given by the observation model. The Shannon entropy $H(\cdot)$ measures the amount of uncertainty about the object model

o and can be written as $H(o_t|f_{1:t}, a_{2:t}) = -\sum_{i=1}^{KQ} P(o^i|f_{1:t}, a_{2:t}) \log P(o^i|f_{1:t}, a_{2:t})$ in the discrete case.

Since (5) is not possible to compute in closed form, we approximate by sampling from the feature distribution, and obtain

$$E[\Delta H(o^i, f_{1:t}, f, a_{2:t})] \approx H(o_t|f_{1:t}, a_{2:t}) - \sum^{\Gamma} H(o_t|f_{1:t}, f, a_{2:t}), \quad (6)$$

with Γ being the number of samples drawn from the feature distribution $P(f|o^i)$. We can further approximate (6) by sampling with zero variance, which yields the mean μ_{o^i} of our trained object model o^i , and with that, a *mean feature* vector $f_{\mu_{o^i}}$. In the case of selecting the next feature, the mean feature vector contains only one feature $f_{\mu_{o^i}}$, which leads to

$$E[\Delta H(o^i, f_{1:t}, f, a_{2:t})] \approx H(o|f_{1:t}, a_{2:t}) - P(f_{\mu_{o^i}}|o^i) H(o|f_{1:t}, f_{\mu_{o^i}}, a_{2:t}). \quad (7)$$

The evaluation of $s_t(a, f)$ is quite costly to compute, since it involves the sequential update of the posterior distribution given the predicted observation. Even with the mean feature approximation of (7), this still involves the computation of KQ predictive posteriors. Restricting attention to only the most probable o^i object hypothesis can further reduce the computational cost.

Stopping Criterion: For LE, we stop integrating more features when the loss of entropy for the selected action is small, $\max(s_t(a, f)) < \tau_1$.

4.1.2 Feature Selection by Mutual Information

In contrast to LE, which uses the predictive posterior distribution, MI uses the conditional entropy. Compared to the predictive posterior distribution, the conditional entropy is much cheaper to evaluate, making the next action fast to compute in comparison. The utility score is computed as $s_t(a, f) = I(o_t; f|a)$, by letting $a = \cdot$, from the averaged MI,

$$I(o_t; f) = \sum_{i=1}^{KQ} P(o^i|f_{1:t}, a_{2:t}) \int p(f|o^i) \log \frac{p(f|o^i)}{\sum_{j=1}^{KQ} p(f|o^j)} df. \quad (8)$$

Since this quantity again is not possible to compute in closed form, we also approximate (8) with the mean feature $f_{\mu_{o^i}}$ by

$$I(o_t; f) \approx \sum_{i=1}^{KQ} P(o^i|f_{1:t}, a_{2:t}) P(f_{\mu_{o^i}}|o^i) \log \frac{P(f_{\mu_{o^i}}|o^i)}{\sum_{j=1}^{KQ} P(f_{\mu_{o^i}}|o^j)}. \quad (9)$$

Here, the mean feature once again is one particular feature, namely the feature we want to test for its information content.

Stopping Criterion: We stop to select a new feature from the current location when $s_t(a, f)$ decreases compared to the previous value [21].

4.2 View Planning

The goal in view planning is to evaluate and select new viewpoints that yield us novel and useful information about the object. Initially, we have not observed anything. So every observation has equal usefulness under the assumption that we do not have any prior knowledge and a uniform belief state. However, every subsequent observation from a different viewpoint bears additional information that may help us to identify the object correctly. The problem consists of not only finding new observation locations that provide previously unseen information, but also viewing positions that help us to improve our current belief of the object in question. In the following subsections, we formulate three methods for viewpoint selection based on three different information-theoretic measures, namely LE, MI and the KL-based *Jeffrey's divergence*.

Stopping Criterion: The stopping criterion for whether to select a new observation location or to terminate depends on the entropy. We will return the MAP estimate of the posterior distribution whenever $H(P(o_t|f_{1:t}, a_{2:t})) < \tau_2$.

4.2.1 Viewpoint Selection by Expected Loss of Entropy

Similar as with LE for feature selection, we compute the next best viewpoint by using the approximation via the mean feature vector $\mathbf{f}_{\mu_{o^i}}$. This time, the mean feature vector contains all features and not just a subset. Since we sample the mean feature from our model, there is no overhead in feature computation, as there would be for making an actual observation. Given the set of actions to move to possible viewing positions, we now compute the utility score $s_t(a, f)$ for an action $a \neq \emptyset$ and a feature f as

$$s_t(a, f) = \sum_{i=1}^{KQ} P(o^i|f_{1:t}, a_{1:t}) E[\Delta H(o^i, f_{1:t}, f, a_{2:t}, a)]. \quad (10)$$

We can see that we again have to compute a complete update of the posterior distribution—this time, however, also including the transition of the sensor from the currently believed state to a new state by performing action a . $E[\Delta H(o^i, f_{1:t}, f, a_{1:t}, a)]$ in (10) is then obtained by a “mean feature” approximation,

$$E[\Delta H(o^i, f_{1:t}, f, a_{2:t}, a)] \approx H(o_t|f_{1:t}, a_{2:t}) - P(\mathbf{f}_{\mu_{o^i}}|o^i, a)H(o_t|f_{1:t}, \mathbf{f}_{\mu_{o^i}}, a_{2:t}, a). \quad (11)$$

4.2.2 Viewpoint Selection by Mutual Information

As seen for feature selection before, the difference between MI and LE is that the second term in (11) is the conditional entropy instead of the entropy of the posterior. The utility score again is directly obtained from MI, $s_t(a, f) = I(o_t; f|a)$, where MI for viewpoint selection is now defined as follows

$$I(o_t; f|a) \approx \sum_{i=1}^{KQ} P(o^i|f_{1:t}, a_{2:t}) P(\mathbf{f}_{\mu_{o^i}}|o^i, a) \log \frac{P(\mathbf{f}_{\mu_{o^i}}|o^i, a)}{\sum_{j=1}^{KQ} P(\mathbf{f}_{\mu_{o^j}}|o^j, a)}. \quad (12)$$

4.2.3 Viewpoint Selection by Jeffrey’s Divergence

Another way to choose new viewing positions is by selecting viewpoints such that, regardless of the true object and pose under observation, the most likely distributions of the resulting measurements are predictably dissimilar. Therefore, a new viewpoint must be selected such that competing hypotheses will appear as different as possible, and that the distinction measured by the features is increased. Finally, the best action is evaluated by computing

$$s_t(a, f) = \sum_{i=1}^{KQ} \sum_{j=1}^{KQ} P(o^i | f_{1:t}, a_{2:t}) P(o^j | f_{1:t}, a_{2:t}) J(P(\mathbf{f}_{\mu_{o^i}} | o^i, a) \| P(\mathbf{f}_{\mu_{o^j}} | o^j, a)), \quad (13)$$

where $J(P(\mathbf{f}_{\mu_{o^i}} | o^i, a) \| P(\mathbf{f}_{\mu_{o^j}} | o^j, a))$ is the Jeffrey’s divergence, defined as $J(P_1 \| P_2) = KL(P_1 \| P_2) + KL(P_2 \| P_1)$, and $KL(P_1 \| P_2) = \int_{-\infty}^{\infty} p_1(x) \log \frac{p_1(x)}{p_2(x)} dx$ denotes the relative entropy or KL-divergence for distributions P_1 and P_2 . The KL-divergence presents the difference in information if samples from the distribution P_1 are assumed to be samples from P_2 .

5 Experiments

To test the framework we have described in the previous sections, we ran extensive evaluations on real world data. In this section, we report on the performance of the individual modules and then demonstrate the performance of the complete system.

5.1 RGB-D Object Dataset

We evaluated our object recognition and action selection framework on the publicly available RGB-D object dataset introduced in [12]. The dataset contains 41 877 data points, each comprising a color image, a 3D point cloud and the orientation of the object, for a total of 300 everyday objects. For each object, the data is recorded from viewing directions emanating from the circumference all around the object, including the three different viewing angles (angles of inclination) of 30°, 45° and 60° per direction. The overall dataset is challenging, since the objects in the dataset can look very similar, especially when they belong to the same object category. Furthermore, the objects in the dataset exhibit large changes in illumination due to the different viewing angles from which the objects were captured.

We compute a feature vector \mathbf{f} that contains $N = 4$ independent features: object bounding box, color histogram, SIFT [14] and viewpoint feature histogram (VFH) [19]. Each feature represents a different object property: size, color, scale-invariant image description and geometry. In accordance with (3), each single feature or component of a feature histogram, respectively, is expressed as independent normal distribution. As we will see, the color feature is the fastest and most powerful single feature among the four features, and thus the first feature we evaluate at each location.

In all our experiments, we use the same evaluation procedure as proposed in [12] and also used in [5]. We test our framework on the task of recognizing the instance

of an object. We train our object models from \mathbf{f} using data captured from viewpoints at viewing angles of 30° and 60° in each direction, and test them against features computed from data that is recorded from viewpoints at viewing angles of 45° . Based on the computed features, we generate $L = 8$ different object models per object instance, which clusters the viewing directions along the circumference horizontally into bins spanning 45° . This subdivision trades computation cost off against performance, and defines the accuracy of the objects' pose estimates. All the reported accuracies of our multi-view object recognition results are expressed as *recall* and are average values over 10 trials. Regarding the stopping criteria, we terminate the action selection procedure whenever either the entropy of the posterior distribution has been reduced to $\tau_2 = 0.1$ or the MAP estimate of the posterior distribution exceeds 60% in accuracy. After termination, we use the MAP estimate of the posterior distribution after (2) to determine the most likely object. For the LE formulation in feature selection we set $\tau_1 = 0.1$ and after termination take as well the MAP estimate as the most likely object.

5.2 Online Feature Selection

The first experiment confirms the usefulness of online feature selection in a single-view recognition task. Selecting the right subset of features can reduce the computation cost while keeping comparable recognition performance.

Each feature has a cost associated for computation; the average computation times of our features are as follows: 0.01 s for the bounding box, 0.0005 s for the color histogram, 0.025 s for the SIFT features, and 2.0 s for the VFH. Fig. 3 shows the average total cost of computation per object and the recognition accuracy plotted against individual feature combinations. The color feature seems to have the best single recognition performance and is the cheapest to compute. The VFH on the other hand does not have a lot of recognition power as such and is rather expensive. From a cost vs. performance standpoint, the feature combination {Color, Bbox, SIFT} performs best with an accuracy of 68.91% and a low average feature computation cost of 0.04 s. However, note that the VFH also has its value since, if included, it increases the recognition accuracy to 74.59%, as detailed in Tab. 1. Additionally, we can see that using feature selection (MI, LE) reduces the cost drastically while keeping recognition accuracy similar to using all features.

Tab. 1 further shows that the average total cost of feature computation can be considerably reduced by making use of feature selection (see the third and fourth row of column "F Cost" in Tab. 1). Here, we utilize LE and MI according to Subsection 4.1 in the feature selection methods. Among the different methods, the LE measure seems to select more features, and more importantly, is much more expensive to compute than MI because we need to update the posterior distribution and compute a full update step.

We compared all results with the results of the original dataset in [12] and [5]. The first paper has very comparable results to our HMM framework and uses similar features to ours, no timing has been provided but it is probably slightly higher than our features. The latter paper reports much higher accuracies; this is mainly due to a more descriptive feature which seems to generalize very well. The cost for such a

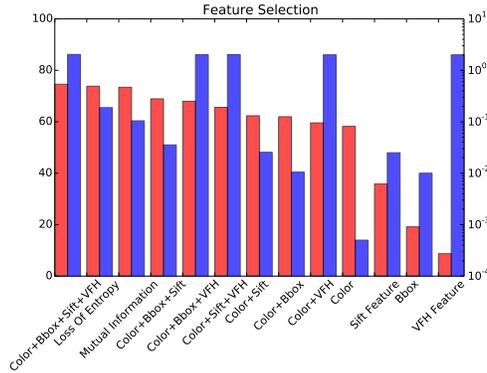


Fig. 3 Performance and cost of different feature combinations in single-view object recognition. The recognition accuracies are displayed as red bars (left side scale, in %). The average cost of computing the feature combinations per object is shown in log scale as blue bars (in seconds).

Method	Accuracy	F Cost (s)	FS Cost (s)	Total (s)
C+B+S	68.91%	0.04	0.00	0.04
C+B+S+V	74.59%	2.04	0.00	2.04
Loss Entropy	73.82%	0.19	1.05	1.24
Mutual Information	73.21%	0.1	0.2	0.3
[12] (kSVM)	74.80%	-	-	-
[4] (Kernel Descriptor)	91.2%	3.2	-	3.2

Table 1 Individual performance for single-view object recognition. “Loss Entropy” and “Mutual Information” denote the feature selection method, with “F Cost” being the average cost per object for computing the features and “FS Cost” being the average cost per object for feature selection.

feature, however, is very high computation time. Moreover, the features used in the two papers are learned in batch, which makes them less flexible compared to our sequential feature learning approach.

5.3 Multi-View Object Recognition

Although the computation time can be substantially reduced through feature selection, the achievable recognition accuracies are limited around 75% for single-view object recognition. The reasons are the relative simplicity of our four feature types and their limitation in generalizing well for the given RGB-D dataset. However, even with fairly simple features we can achieve improved recognition accuracies for both object class and pose by utilizing an information-theoretic multi-view object recognition approach. In the following tests, we use the MI measure for feature selection, since it is on par with the LE measure regarding recognition performance, and in addition, has lower computation cost (see Tab. 1).

Tab. 2 shows the performance and Tab. 3 summarizes the cost of multi-view object recognition. The first row shows that using all features (i.e., no feature selection) has a very high cost but random viewpoint selection already increases the recogni-

Method	Accuracy	Observations	Avg. Bin Error
C+B+S+V (no FS)	84.48% \pm 1.4	3.4 \pm 0.2	2.41 \pm 0.07 (bins) (\pm 67.5 $^\circ$)
Random	86.76% \pm 0.4	3.2 \pm 0.2	2.19 \pm 0.02 (bins) (\pm 67.5 $^\circ$)
Jeffrey’s Divergence	88.34% \pm 0.6	2.7 \pm 0.1	1.62 \pm 0.01 (bins) (\pm 45.0 $^\circ$)
Loss Entropy	89.82% \pm 1.3	2.4 \pm 0.2	1.54 \pm 0.01 (bins) (\pm 45.0 $^\circ$)
Mutual Information	91.87% \pm 0.7	2.5 \pm 0.1	1.48 \pm 0.01 (bins) (\pm 45.0 $^\circ$)

Table 2 Performance of our multi-view object recognition framework. All but the first row uses the MI measure in the feature selection method. We show the average recognition accuracy, the average number of observations and the average bin error of estimated viewing directions. Each bin represents a range of observation angles of 45 $^\circ$.

Method	F Cost (s)	FS Cost (s)	VS Cost (s)	Total Cost (s)
C+B+S+V (no FS)	7.02 \pm 0.42	0.00	0.0	7.02 \pm 0.42
Random	0.40 \pm 0.04	0.83 \pm 0.05	0.0	1.23 \pm 0.09
Jeffrey’s Divergence	0.30 \pm 0.03	0.58 \pm 0.02	3.E-4 \pm 2.E-4	0.88 \pm 0.05
Loss Entropy	0.23 \pm 0.04	0.50 \pm 0.04	5.4 \pm 0.6	6.13 \pm 0.68
Mutual Information	0.24 \pm 0.03	0.51 \pm 0.03	1.1 \pm 0.2	1.85 \pm 0.26

Table 3 Average computation cost per object for our multi-view object recognition framework. The columns show the average cost for computing the feature (“F Cost”), average cost for feature selection (“FS Cost”) and average cost for viewpoint selection (“VS Cost”).

tion accuracy by about 10%. In the second row, due to feature selection using mutual information and random viewpoint selection, we can see a gain in reduced computation time and an increase in accuracy. On average, random viewpoint selection terminates after a little more than 3 observations. By using the information-theoretic measures for viewpoint selection, we find that LE, MI and Jeffrey’s divergence all perform similarly well—with MI being the best in terms of recognition accuracy. The lowest computation time is achieved by using the Jeffrey’s divergence because it can compute new viewing positions in almost no time. The MI formulation is superior to the LE formulation in speed, since it does not need to compute the predictive posterior.

We further evaluate our multi-view approach for estimating an object’s pose. As previously mentioned, we train each object with a number of different object models, in our case $L = 8$ models per object. When observing a feature, we can infer from which model it has most likely been generated, and by this, further infer about the orientation of the object. The performance of the object’s pose estimation is presented in Tab. 2. The table shows the average bin error of the estimated viewing directions, which is computed from the differences between the inferred and the true viewing angles. For the first row, all the features were computed (i.e., no feature selection); the feature selection methods with different information-theoretic measures were used for the remaining rows. We are on average > 2 bins off in our estimation without view planning (i.e., random selection of views), which corresponds to an average pose error of about 67.5 $^\circ$. If we apply one of the three more sophisticated viewpoint selection methods, we yield an average bin error of < 2 bins. This results



Fig. 4 Quadcopter experiment and test setup. The quadcopter robot, equipped with a RGB-D camera, collects multiple views of a target object (e.g., a cereal box) in an active object recognition task.

in an average pose error of 45° . For comparison, the experiments in [5] report an average pose error of 44.8° , which is on par with our results.

5.4 Quadcopter Experiments

Finally, to show the overall performance of our active multi-view object recognition framework, we have conducted experiments on our custom-built quadcopter platform. As shown in Fig. 4, the quadcopter is equipped with an Xtion RGB-D camera for data acquisition and can fly autonomously to a new observation location.

Given an object, the task is to fly the quadcopter around the object and take as many observations as necessary to correctly recognize the object. Compared to the simulated experiments above, where data is captured with a RGB-D camera mounted on a tripod, experiments on a robotic platform are generally more challenging. Specific challenges regarding object recognition on a fast moving platform include motion blur and constant changes in illumination.

As we do not have any of the objects used in the RGB-D dataset at hand, we captured 10 new objects and added them to the dataset, which now consists of a total of 310 objects. We captured the data in the same way as described in [12] and chose objects similar to the ones already contained in the dataset. We trained our models with all 310 objects and observations taken at the two viewing angles of 30° and 60° . In the actual experiment, the quadcopter then flies around the object at a distance and height that allows for a viewing angle of 45° roughly.

In each run of the experiments, the robot starts at a random position with respect to the object, from where it acquires the first observation of the object. From there on, it selects the following actions as described in the previous section—either inferring about another feature from \mathbf{f} or flying to a new viewpoint. We terminate the process once our stopping criteria are reached. Fig. 4 shows an example of the quadcopter taking observations of a cereal box object during one of the experiments.

Tab. 4 summarizes the experimental results from the quadcopter experiments, averaged over eight runs with different starting locations. We can see that for all experiments, whether conducted with single-view or multi-view object recognition, online feature selection not only helps to decrease the computation time but also increases the recognition accuracy. The average pose error is similar to the error in the simulation experiments, however, the average bin error is smaller, which results in more correctly classified bins in total.

Multi-view recognition on average only required 1.6 observations, which means that often we are able to make a confident decision after only one scan. In some

Method	Accuracy	Obsv.	Avg. Bin Error	Avg. Cost (s)
Single-View				
C+B+S+V (no FS)	82.16%	1	2.1(bins)($\pm 67.5^\circ$)	2.14
Feature Selection (MI)	87.32%	1	1.9(bins)($\pm 45^\circ$)	0.19
Multi-View				
Mutual Information (no FS)	96.5% \pm 0.5	2.0 \pm 0.1	1.6 \pm 0.1(bins)($\pm 45^\circ$)	5.63 \pm 0.3
Mutual Information	99.8% \pm 0.1	1.6 \pm 0.01	1.1 \pm 0.1(bins)($\pm 45^\circ$)	1.48 \pm 0.2

Table 4 Single-view and multi-view object recognition based on data from the quadcopter robot. In both cases, results with and without feature selection are presented. Advantages of online feature selection as well as multi-view object recognition are revealed by the experiments.

instances, however, by acquiring additional observations, the recognition accuracy could be increased by more than 10% to 99.8%. This of course comes with a price of actually moving the robot to a new location which furthermore increases the cost of recognizing the object. In practice, whether the cost of moving is justified, is probably a tradeoff between accuracy and energy saving. However, we can clearly see that adaptive feature selection has huge benefits by decreasing computation time, which is very beneficial for vehicles with low energy budget like quadcopters.

6 Conclusion and Future Work

We have presented an information-theoretic action selection framework for object recognition that is capable of sequential training and classification, view planning as well as online feature selection. Online feature selection is especially useful since it is adaptive, can reduce computation time drastically and is more flexible, compared to offline feature selection. Our framework was shown to lead to increased recognition accuracies by using four relatively simple features in combination with the multi-view object recognition approach. Furthermore, we demonstrated that information-theoretic formulations perform better than random action selection strategies in that they require fewer observations and less computation time.

Compared to state-of-the-art single-view object recognition methods, our multi-view object recognition approach achieves comparable results. Our approach is not primarily meant to replace single-view methods but rather to complement them. Put in other words, if a confident decision can be made about an object after only one observation, there will be no need for taking a second view. However, there will always be objects with inherent ambiguity, and it is in those cases in which the multi-view approach is brought to bear. Feature selection is useful regardless of the number of views and the object in question, since it reduces the dimensionality of the feature space and decreases the overall computation cost.

The kernel features used in [5] seem to be far superior to our simple features. Thus, the question arises how our approach would perform when using these more complex features instead. We are currently in the process of experimenting with kernel features and will test them with our framework.

Acknowledgements: The authors would like to thank Jörg Müller for his help with the development of the quadcopter robot platform.

References

1. Ali, H., Marton, Z.C.: Evaluation of feature selection and model training strategies for object category recognition. In: *Intelligent Robots and Systems*, pp. 5036–5042. IEEE/RSJ (2014)
2. Atanasov, N., Sankaran, B., Ny, J.L., Pappas, G.J., Daniilidis, K.: Nonmyopic View Planning for Active Object Classification and Pose Estimation. *IEEE Transactions on Robotics* **30**(5), 1078 – 1090 (2014)
3. Bajcsy, R.: Active perception. *Proceedings of the IEEE* **76**(8), 966–1005 (1988)
4. Bo, L., Ren, X., Fox, D.: Depth kernel descriptors for object recognition. In: *Int. Conf. on Intelligent Robots and Systems*, pp. 821–826. IEEE (2011)
5. Bo, L., Ren, X., Fox, D.: Unsupervised feature learning for rgb-d based object recognition. In: *Experimental Robotics - The 13th Int. Symp. on Exp. Robotics*, pp. 387–402 (2012)
6. Borotschnig, H., Paletta, L., Prantl, M., Pinz, A.: Active Object Recognition in Parametric Eigenspace. In: *Proc. BMVC*, pp. 63.1–63.10 (1998)
7. Brown, M., Lowe, D.G.: Unsupervised 3D object recognition and reconstruction in unordered datasets. In: *5th Int. Conf. on Digital Imaging and Modeling*, pp. 56–63 (2005)
8. Denzler, J., Brown, C.M.: Information Theoretic Sensor Data Selection for Active Object Recognition and State Estimation. *Trans. Pattern Anal. Mach. Intell.* **24**(2), 145–157 (2002)
9. Finman, R., Whelan, T., Kaess, M., Leonard, J.J.: Toward lifelong object segmentation from change detection in dense RGB-D maps. In: *European Conf. on Mobile Robots* (2013)
10. Grabner, H., Bischof, H.: On-line Boosting and Vision. In: *Conf. on Computer Vision and Pattern Recognition*, vol. 1, pp. 260–267 (2006)
11. Hollinger, G.A., Mitra, U., Sukhatme, G.S.: Active Classification: Theory and Application to Underwater Inspection. In: *Int. Symposium on Robotics Research (ISRR)*. Flagstaff, AZ, (2011)
12. Lai, K., Bo, L., Ren, X., Fox, D.: A large-scale hierarchical multi-view RGB-D object dataset. In: *Int. Conf. on Robotics and Automation*, pp. 1817–1824. IEEE (2011)
13. Laporte, C., Arbel, T.: Efficient discriminant viewpoint selection for active Bayesian recognition. *International Journal of Computer Vision* **68**(3), 267–287 (2006)
14. Lowe, D.G.: Distinctive Image Features from Scale-Invariant Keypoints. *Int. J. Comput. Vision* **60**(2), 91–110 (2004)
15. Luber, M., Spinello, L., Arras, K.O.: People Tracking in RGB-D Data with On-line Boosted Target Models. In: *Int. Conf. on Intelligent Robots and Systems*, pp. 3844–3849. IEEE/RSJ (2011)
16. Paletta, L., Pinz, A.: Active object recognition by view integration and reinforcement learning. *Robotics and Autonomous Systems* **31**(1), 71–86 (2000)
17. Peng, H., Long, F., Ding, C.: Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy. *Transactions on Pattern Analysis and Machine Intelligence* **27**(8), 1226–1238 (2005)
18. Potthast, C., Sukhatme, G.S.: A Probabilistic Framework for Next Best View Estimation in a Cluttered Environment. *J. Vis. Comun. Image Represent.* **25**(1), 148–164 (2014)
19. Rusu, R.B., Bradski, G., Thibaux, R., Hsu, J.: Fast 3D Recognition and Pose Using the Viewpoint Feature Histogram. In: *Int. Conf. on Intelligent Robots and Systems*. IEEE/RSJ (2010)
20. Salvagnini, P., Pernici, F., Cristani, M., Lisanti, G., Del Bimbo, A., Murino, V.: Non-myopic information theoretic sensor management of a single pan-tilt-zoom camera for multiple object detection and tracking. *Computer Vision and Image Understanding* **134**, 74–88 (2015)
21. Verleysen, M., Rossi, F.: *Advances in Feature Selection with Mutual Information*. CoRR (2009)
22. Zhou, X.S., Comaniciu, D., Krishnan, A.: Conditional feature sensitivity: a unifying view on active recognition and feature selection. *9th Int. Conf. on Computer Vision* (2003)