

# Multi-Armed Bandit Models for Sample-Based Grasp Planning in the Presence of Uncertainty [v-14 09-30-2014]

Michael Laskey<sup>1</sup>, Zoe McCarthy<sup>1</sup>, Jeff Mahler<sup>1</sup>, Florian T. Pokorny<sup>3</sup>, Sachin Patil<sup>1</sup>,  
Jur Van Den Berg<sup>4</sup>, Danica Kragic<sup>3</sup>, Pieter Abbeel<sup>1</sup>, Ken Goldberg<sup>2</sup>

**Abstract**—In unknown or changing environments such as warehouses for consumer goods, robots may need to quickly plan grasps on previously unknown objects. Due to sensor noise, occlusions, and surface properties such as transparency and specularly, it may be difficult to plan successful grasps. Planning grasps in the presence of uncertainty can be computationally demanding because fully evaluating the expected quality of even a single grasp may require a large number of samples. In this paper, we show that framing grasp planning as a multi-armed bandit problem, a standard framework for making sequential decisions, provides a principled way to quickly select a grasp from a set of candidates on an uncertain shape. We demonstrate this method for planning grasps in the presence of shape uncertainty encoded by a Gaussian process implicit surface (GPIS), which can be especially computationally demanding. We utilize successive elimination to solve a multi-armed bandit problem and show a speedup of T versus planning grasps using existing methods.

## I. INTRODUCTION

Consider a robot packing boxes in a shipping warehouse environment, where it may frequently encounter new consumer products and need to process them quickly. The robot may need to rapidly plan grasps for these objects without prior knowledge of their shape and pose. However, the robot may not be able to measure these quantities exactly due to sensor noise and missing data due to partial visibility and object properties such as transparency. Grasp planners that assume an exact measurement of shape, pose, or its own state may fail.

This motivates using knowledge of uncertainty to select grasps, but most current methods for evaluating the quality of a single grasp in the presence of uncertainty require use an exhaustive sampling over the possible values of the uncertain quantity [18], [19], [40]. To select a grasp with high quality this evaluation is often performed for a large set of potential grasps, which can be very time-consuming. However, when ranking a set of grasps we may be able to determine the relative quality between grasps with only a few samples and

throw away grasps that are likely to be suboptimal [17]. Thus, we can adaptively concentrate grasp evaluation on the grasps that are most likely to have the highest quality based on the evaluation done so far.

The multi-armed bandit (MAB) problem, a framework for sequential decision making problems, provide a principled way to reason about selecting the next grasp to evaluate and the grasps to discard from consideration [5], [21], [32]. The goal in MAB problems is to make a sequence of decisions over a set of possible options such that a measure of the expected value of such decisions is maximized. Solutions to the MAB problem are particularly useful in applications where it is too expensive to fully evaluate a set of options; for example, in optimal design of clinical trials [34], market pricing [33], and choosing strategies for games [36].

[TODO: REF FIGURE]

Our main contribution in this paper is formulating the problem of planning grasps with a high expected quality in the presence of uncertainty as a multi-armed bandit problem. We use this formulation to rank a set of potential grasps by expected Ferrari-Canny quality under shape uncertainty represented as a Gaussian process implicit surface (GPIS), a Bayesian representation of shape uncertainty that seen recent use in various robotic applications [11], [14] but has been limited by its computational complexity [30], [41]. We also show how to estimate distributions on the contact points, surface normals, and center of mass using a GPIS shape representation. Our experiments demonstrate that using the MAB sampling method improves the time to rank a set of N grasps by T over other grasp evaluation methods.

## II. RELATED WORK

The multi-armed bandit (MAB) problem [5], [21], [32] considers maximize the total expected reward over a sequence of possible choices over a set of competing options, each of which has a probabilistic reward function. Solutions to this problem, called *policies*, include choosing the next option based on Thompson sampling [2], the Gittins index policy [39], and by using an upper confidence bound (UCB) for the expected reward of each option [4]. The UCB can be computed from the sample mean [1], Kullback-Leibler divergence [9], or a prior on the rewards [16]. Solutions to the multi-armed bandit problem have been used in applications for which evaluating all possible options is expensive or impossible, such as the optimal design of clinical trials [34], market pricing [33], and choosing strategies for games [36]. In this paper we consider the "budgeted multi-arm bandit"

<sup>1</sup>Department of Electrical Engineering and Computer Sciences; {mdlaskey, zmccarthy, jmahler, sachinpatil, pabbeel}@berkeley.edu

<sup>2</sup>Department of Industrial Engineering and Operations Research and Department of Electrical Engineering and Computer Sciences; goldberg@berkeley.edu

<sup>1-2</sup> University of California, Berkeley; Berkeley, CA 94720, USA

<sup>3</sup>Computer Vision and Active Perception Lab, Centre for Autonomous Systems, School of Computer Science and Communication, KTH Royal Institute of Technology, Stockholm, Sweden {fpokorny, dani}@kth.se

<sup>4</sup>Google; Amphitheatre Parkway, Mountain View, CA 94043, USA jurvandenber@gmail.com

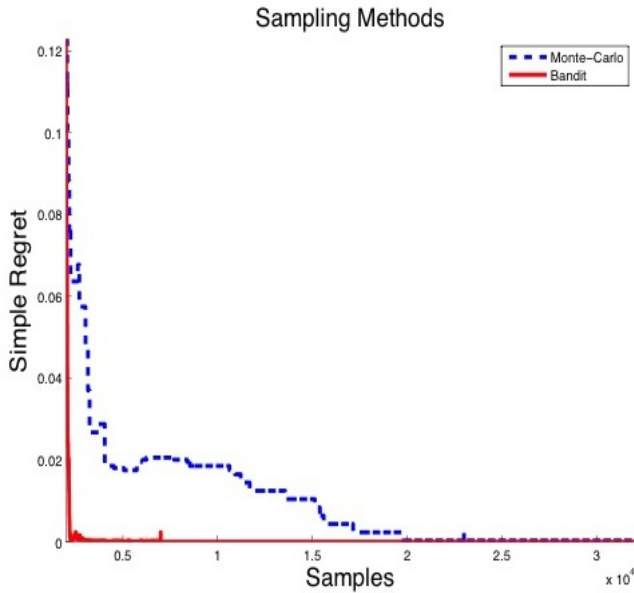


Fig. 1: Convergence in regret of the bandit sampling method (red), compared to the traditional Monte-Carlo method. The fast convergence of the bandit method is due to its ability to intelligently pick what grasp to sample next in a given set of proposed grasps on an object with shape uncertainty

problem, where exploration and exploitation are decoupled and the agent can only explore for a finite time before deciding which arm to exploit [24].

Past work on grasping under uncertainty has considered state uncertainty [13], [37], uncertainty in contact locations with an object [42], uncertainty in object pose [10], [40], [19]. The effect of uncertainty in object geometry on grasp selection has been studied for spline representations of objects [10], extruded polygonal mesh models [17], [18], and point clouds [15].

Currently, the most common method of evaluating grasp quality under shape and pose uncertainty is to rank a set of random grasps on an object using Monte-Carlo integration on shape and pose samples to evaluate a quality measure [10], [17], [18]. Monte-Carlo integration involves drawing random samples from a distribution to approximate an integral [8], which can be slow when the distribution is high-dimensional, such as for distributions on possible shape. To address this, Kehoe et al. [17] demonstrated a procedure for finding a minimum bound on grasp quality given shape uncertainty, which reduced the number of terms needed in Monte-Carlo integration in order to rank grasps, but the adaptive sampling method involves several hyperparameters whereas our proposed budgeted MAB solution does not. Laaksonen et al. [20] used Markov Chain Monte Carlo (MCMC) sampling to estimate grasp stability and object pose online under shape and pose uncertainty. MCMC simplified sampling from complicated distributions on pose and shape, but it can be slow to converge to the correct distribution [3].

We study our MAB sampling method using Gaussian process implicit surfaces (GPIS) to represent shape uncertainty, based on the ability to combine various modes of

noise observations such as tactile, laser and visual [30], [41], [11] and its recent use in modeling uncertainty for a number of robotic applications. Hollinger et al. used GPIS to perform active sensing on the hulls in underwater boats [14]. Dragiev et al. showed how GPIS can represent shapes for grasps and used as a grasp controller on the continuous signed distance function [11]. Mahler et al. used GPIS representation to find locally optimal anti-podal grasps by framing grasp planning as an optimization problem [25]. However, sampling from GPIS for grasp quality evaluation is computationally intensive because it involves the Cholesky decomposition of an  $m \times m$  matrix which takes  $O(m^6)$  time, where  $m$  is the dimension for a 2D grid over an object [27]. Therefore, GPIS is a good illustration of benefits provided by our efficient MAB approach to grasp planning.

### III. MULTI-ARM BANDIT PROBLEM

The multi-arm bandit problem originally described by Robbins [?] is a statistical decision model of an agent trying to make correct decisions, while gathering information at the same time. The traditional setting of a multi-arm bandit problem is a gambler that has  $K$  independent slot machine arms and tries to infer which one will yield the highest reward. A successful gambler would want to exploit the machine that currently yields the highest reward and explore new arms to see if they give better rewards. Developing a policy that successfully trades between exploration and exploitation has been the focus of extensive research, since the problem formulation [7], [32], [6].

There are a number of algorithms for developing policies to balance exploration and exploitation. One algorithm is  $\epsilon$ -greedy, which is the idea of choosing the arm with the highest empirical expected reward with  $1 - \epsilon$  probability and choosing a random arm with probability  $\epsilon$  [5]. A class of algorithms that have theoretical guarantees are Upper Confidence Bound (UCB). UCB algorithms maintain the empirical expected reward based off of pulling each arm multiple times, while also estimating an upper bound on the true expected reward using assumptions on the probability distribution of rewards and number of times has been sampled for each arm. The algorithm chooses the arm at each step which has the highest upper bound on expected reward, since that is the most potentially promising arm. When the rewards come from exponential family distribution, UCB minimizes cumulative regret, which is the number of times a sub-optimal arm is pulled times the difference between the true expected reward of an optimal arm and the true expected reward of that sub-optimal arm. In practice, when the distribution on the rewards of arms is not known, the empirical methods such as  $\epsilon$ -greedy have shown to have better performance in some situations [?].

We can formulate our problem as a "budgeted multi-arm bandit problem" [24]. The budgeted multi-arm bandit problem has a finite stopping time. The objective of the budgeted multi-arm bandit problem is to minimize the "simple regret", which is the difference between the true expected value of the actual best arm and the true expected value of the arm

pulled at the stopping time. Thus the budgeted multi-arm bandit problem reduces to exploration until the very last time step, at which time one exploitation step is taken.

#### IV. PROBLEM DEFINITION

Before we present the problem definition, we introduce the grasping model, the line of action.

##### A. Line of action

Similar to the work of [10], we assume that each gripper finger approaches along a *line of action*, a 1D curve  $\gamma(t)$  with endpoints  $a$  and  $b$  as seen in Fig. 2. A gripper finger starts at point  $a$  and moves towards  $b$ , we assume  $a$  is far enough away to be collision free of the object. Each gripper contact is defined by a line of action, so we assume the following tuple is provided  $\Gamma = (\gamma_1(\cdot), \dots, \gamma_m(\cdot))$ , which designates a proposed *grasp plan*. [TODO: MOVE THE NEXT TWO LINES] While we currently assume the gripper moves free of noise, this approach would be applicable to high precision robots such as industrial based on. Future work will look at how to efficiently sample when the approach trajectory has noise.

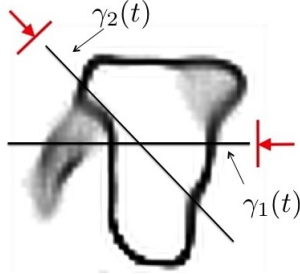


Fig. 2: Illustration of a grasp plan  $\Gamma$  composed of two lines of action,  $\gamma_1(t)$  and  $\gamma_2(t)$

##### B. Problem Definition

Given a 2-D workspace  $\mathcal{W}$ , with an unknown object represented as a trained GPIS model (we will describe the GPIS model later) and set of possible grasp plans  $G$ . We are interested in determining

$$\Gamma^* \in \operatorname{argmax}_{\Gamma \in G} E(Q(\Gamma)) \quad (1)$$

with respect to a chosen grasp metric  $Q$ .

#### V. MULTI-ARM BANDITS FOR GRASP SELECTION

While a standard approach to solving the problem in Eq. 1 would be to simply perform Monte-Carlo integration on each  $\Gamma_i$  and compute the expected grasp quality, we propose treating the problem as a multi-arm bandit problem and forming a policy for selecting which grasp to sample. In our setting, we have a probabilistic shape representation and would like to evaluate many potential grasps on that shape model. Motivated by limited computational resources we are interested in how to intelligently allocate sampling resources to efficiently find the best grasp plan  $\Gamma^*$ . Here each arm corresponds to a different grasp plan and pulling the arm is sampling the shape representation and evaluating the arm's

grasp plan on the sampled shape representation. The reward for pulling an arm is the grasp quality of the resulting grasp on the sampled shape. We have a policy for exploration of different grasp plans (i.e. choosing which one to sample next) and at a given stopping time we choose to execute the grasp plan with the highest expected quality based on the samples received so far, which would correspond to actually performing the grasp. The number of samples needed before the simple regret reaches zero, determines how effective an exploration policy is for grasp evaluation.

In solving Eq. 1, we want to pick the grasp plan with the highest  $E(Q(\Gamma))$  determined via Monte-Carlo Integration. Due to the non-linear nature of our grasp metric, we currently do not have a way to represent the distribution on grasp quality of  $p(Q(\Gamma))$  in closed form. Thus, we are restricted in the type of bandit policies available to us, since most sophisticated policies require distributional information. [TODO: DOUBLE CHECK WITH DYLAN]

For our exploration policy, we maintain a subset of the possible grasp plans that are statistically indistinguishable from the best grasp plan. During exploration, we sample uniformly from this subset. This has been shown to have good performance when the number of evaluations is large [7]. For each grasp plan, we maintain a 95% confidence interval around the expected quality of the grasp plan. The width of the confidence interval for a grasp plan  $i$  that has been sampled  $n_i$  times and the samples have standard deviation  $\sigma_i$  is:

$$C_i := \frac{1.96\sigma_i}{\sqrt{n_i}}. \quad (2)$$

Let  $\hat{\mu}_i$  be the sample mean of the grasp quality for grasp plan  $i$ . Then the confidence interval for grasp plan  $i$  is  $[\hat{\mu}_i - C_i, \hat{\mu}_i + C_i]$  [8]. After a sample from a grasp plan, we check to see if the confidence interval of the sampled grasp intersects with the confidence interval of the current best grasp plan, and if it does not, we prune that grasp plan from the current active set of grasp plans. Since we prune grasps that are below a confidence interval we have statistical guarantees unlike earlier approaches to this problem[18]. In future work, we will look at using other exploration strategies.

#### VI. EVALUATING A GRASP

In order to solve our problem definition, we must evaluate  $E(Q(\Gamma))$  for a given grasp plan  $\Gamma$ . We will first discuss how the GPIS is constructed, then which grasp metric  $Q$  we chose and lastly proceed into evaluating the expectation  $E$  efficiently.

##### A. Gaussian Process (GP) Background

Given the training data  $\mathcal{D} = \{\mathcal{X}, \mathbf{y}\}$  and covariance function  $k(\cdot, \cdot)$ , the posterior density  $p(sd_* | \mathbf{x}_*, \mathcal{D})$  at a test point  $\mathbf{x}_*$  is shown to be [31]:

$$\begin{aligned} p(sd_* | \mathbf{x}_*, \mathcal{D}) &\sim \mathcal{N}(\mu(\mathbf{x}_*), \Sigma(\mathbf{x}_*)) \\ \mu(\mathbf{x}_*) &= k(\mathcal{X}, \mathbf{x}_*)^\top (K + \sigma^2 I)^{-1} \mathbf{y} \\ \Sigma(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathcal{X}, \mathbf{x}_*)^\top (K + \sigma^2 I)^{-1} k(\mathcal{X}, \mathbf{x}_*) \end{aligned}$$

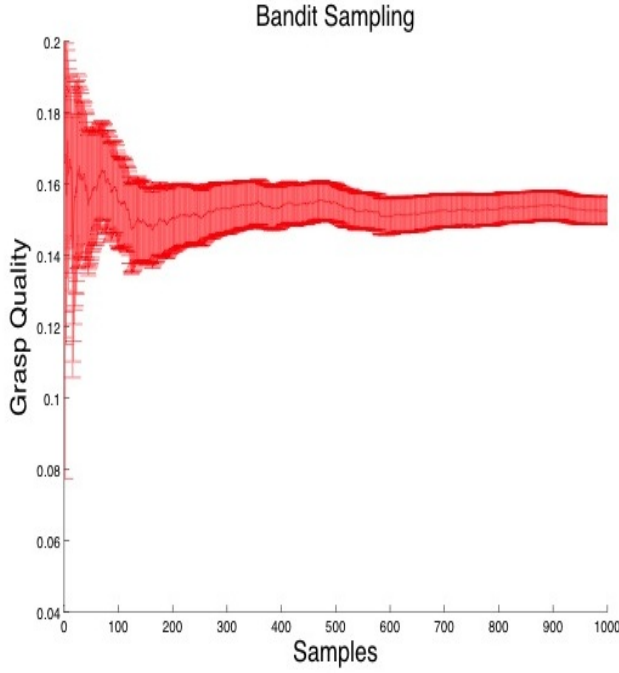


Fig. 3: Convergence of expected grasp quality,  $E(Q(g))$ , for a typical grasp. You can see it can take hundreds of evaluations to correctly estimate the expected grasp quality, thus it is important to intelligently allocate sampling resources

where  $K \in \mathbb{R}^{n \times n}$  is a matrix with entries  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$  and  $k(\mathcal{X}, \mathbf{x}_*) = [k(\mathbf{x}_1, \mathbf{x}_*), \dots, k(\mathbf{x}_n, \mathbf{x}_*)]^\top$ . This derivation can also be used to predict the mean and variance of the function gradient by extending the kernel matrices using the identities [35]:

$$\text{cov}(sd(\mathbf{x}_i), sd(\mathbf{x}_j)) = k(\mathbf{x}_i, \mathbf{x}_j) \quad (3)$$

$$\text{cov}\left(\frac{\partial sd(\mathbf{x}_i)}{\partial x_k}, sd(\mathbf{x}_j)\right) = \frac{\partial}{\partial x_k} k(\mathbf{x}_i, \mathbf{x}_j) \quad (4)$$

$$\text{cov}\left(\frac{\partial sd(\mathbf{x}_i)}{\partial x_k}, \frac{\partial sd(\mathbf{x}_j)}{\partial x_l}\right) = \frac{\partial^2}{\partial x_k \partial x_l} k(\mathbf{x}_i, \mathbf{x}_j) \quad (5)$$

For our kernel choice we decided to use the square exponential kernel, similar to [11]. Other kernels relevant to GPIS are the thin-plate splines kernel and the Matern kernel [41].

We construct a GPIS by learning a Gaussian Process to fit measurements of a signed distance field of an unknown object. Precisely,  $x_i \in \mathbb{R}^2$  in 2D and  $x_i \in \mathbb{R}^3$  in 3D, and  $y_i \in \mathbb{R}$  is a noisy signed distance measurement to the unknown object at  $x_i$ .

### B. Grasp Metric

In their pioneering work over two decades ago Ferrari and Canny [12], demonstrated a method to rank grasps by considering their contact points and surface normals. Importantly the magnitude of  $Q$  yields a measurement that allows one to rank grasps by their physical stability and evaluate the property of force-closure. Furthermore, it has wide spread use in grasp packages like GraspIT[26], OpenGrasp[22] and Simox [38], which motivates studying its effect with uncertainties.

The  $L^1$  version of the metric works by taking as input the contact points  $\mathbf{c}_1, \dots, \mathbf{c}_m$ , surface normals  $\mathbf{n}_1, \dots, \mathbf{n}_m$ , center of mass  $\mathbf{z}$  and friction coefficient  $\mu$ . Then constructing a convex hull around the wrenches made up of those parameters and finding the radius of the largest unit ball centered at the origin in wrench space. A wrench is defined as concatenation of a force and torque vector. If the convex hull doesn't enclose the origin, the grasp is not in force-closure. Thus a grasp can be parameterized by the following tuple  $g = \{\mathbf{c}_1, \dots, \mathbf{c}_m, \mathbf{n}_1, \dots, \mathbf{n}_m, \mu, \mathbf{z}\}$ , our method is applicable to all grasp metrics that represent a grasp as the tuple  $g$ , such as [10], [23].

### C. Calculating the Expected Grasp Quality

Given a proposed grasp plan  $\Gamma$ , the expected grasp quality can be evaluated as follows:

$$E(Q(\Gamma)) = \int Q(g|S, \Gamma) p(S) dS \quad (6)$$

Where  $Q(g|S, \Gamma)$  is the grasp quality that is computed on a shape sample drawn from  $p(S)$ . To compute this we intersect the zero crossing of the level set with the proposed grasp plan  $\Gamma$  and determine the parameters  $g$ , this has been the approach taken in previous work [17], [18], [10]. See Fig. 4 for an example of what samples drawn from  $p(S)$  induced by GPIS look like. For computational reasons we approximate the integral via Monte-Carlo Integration. We use importance sampling to draw from the distribution induced by GPIS and calculate the following

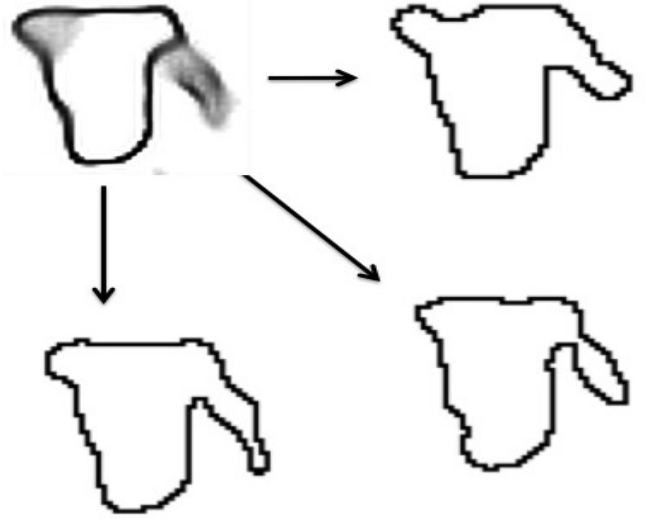


Fig. 4: Shape samples drawn from Eq. 7 on the object in the upper left corner. Given a shape sample we highlight the zero-crossing of the level set in black

$$E(Q(\Gamma)) \approx \frac{1}{N} \sum_{i=1}^N Q(g|S_i, \Gamma), \quad S_i \sim p(S)$$

To compute the above distribution we must draw samples from  $p(S)$ . In order to draw shape samples from a GPIS, one needs to sample from signed distance function,  $sd$ , over the



joint on all points in the workspace  $\mathcal{W}$  or  $p(sd(\mathcal{W}))$ . Since this is a GPIS, we know the following

$$p(S) = p(sd(\mathcal{W})) \sim N(\mu(\mathcal{W}), \Sigma(\mathcal{W})) \quad (7)$$

Thus if the workspace is an  $n \times n$  grid, the joint distribution is an  $n^2$  multi-variate Gaussian, due to  $sd : \mathbb{R}^2 \rightarrow \mathbb{R}$ . Sampling from a Gaussian involves inverting the covariance matrix and inversion is in the naive way  $O(n^3)$  [29]. Thus the complexity of this operation is  $O(n^6)$  in 2D and  $O(n^9)$  in 3D.

To reduce complexity we propose sampling not from the shape distributions, but instead from the distributions on the grasps parameters themselves. We recall that a grasp according to our metric is defined as the tuple  $g = \{\mathbf{c}_1, \dots, \mathbf{c}_m, \mathbf{n}_1, \dots, \mathbf{n}_m, \mu, z\}$ . We are thus interested in calculating  $p(g|\Gamma, \mu(x), \Sigma(x))$ . The distribution on a grasp is defined then as:

$$p(g) = p(\mathbf{c}_1, \dots, \mathbf{c}_m, \mathbf{n}_1, \dots, \mathbf{n}_m | \Gamma, \mu(x), \Sigma(x)) \quad (8)$$

We note here that we currently use the friction coefficient  $\mu$  and the expected center of mass  $\bar{z}$  as deterministic values. For grippers that do not approach along the same line of action (i.e. non-parallel jaw grippers) we make the assumption that each contact and normal pair is independent, or

$$p(g) = \prod_{i=1}^m p(\mathbf{c}_i, \mathbf{n}_i | \gamma_i(t), \mu(x), \Sigma(x)) \quad (9)$$

We compute the expected grasp quality now as follows:

$$E(Q(\Gamma)) = \frac{1}{N} \sum_{i=1}^N Q(g_i), \quad g_i \sim p(g) \quad (10)$$

We will now show how these distributions can be computed and how the complexity for evaluating a grasp is reduced from  $O(n^6)$  to  $O(n^3)$ .

## VII. DISTRIBUTION OF GRASP PARAMETERS

To sample from  $p(g)$ , we need to sample from the distributions associated with a line of action  $p(\mathbf{n}_i, \mathbf{c}_i | \gamma_i(t), \mu(x), \Sigma(x))$ . Using Bayes rule we can rewrite this as

$$p(\mathbf{n}_i, \mathbf{c}_i | \gamma_i(t), \mu(x), \Sigma(x)) = p(\mathbf{n}_i | \mathbf{c}_i, \gamma_i(t), \mu(x), \Sigma(x)) p(\mathbf{c}_i | \gamma_i(t), \mu(x), \Sigma(x))$$

In section VII-A, we look at how to sample from  $p(\mathbf{c}_i | \gamma_i(t), \mu(x), \Sigma(x))$ . Then in section VII-B, we look at how to sample from  $p(\mathbf{n}_i | \mathbf{c}_i, \gamma_i(t), \mu(x), \Sigma(x))$  and present a novel visualization technique for the distribution on surface normals. Lastly in section VII-C, we show a way to calculate the expected center of mass assuming a uniform mass distribution. .

### A. Distribution on Contact Points

We would like to find the distribution on contact point  $\mathbf{c}_i$ . A contact point in terms of the GPIS and line of action model can be defined as the point,  $t$ , when the signed distance function is zero and no points before said point along the line have touched the surface. We express this as the following conditions:

$$sd(\gamma(t)) = 0 \quad (11)$$

$$sd(\gamma(\tau)) > 0, \quad \forall \tau \in [a, t) \quad (12)$$

We will now demonstrate how to efficiently sample from  $p(\mathbf{c}_i | \mu(x), \Sigma(x), \gamma_i t)$

The probability distribution along the line  $\gamma(t)$  is given by:

$$p(sd(\gamma(t)); \mu(t), \Sigma(t)) \quad \forall t \in [a, b] \quad (13)$$

This gives the signed distance function distributions along the entire line of action in the workspace as a multivariate Gaussian. One could think of this as a marginalization of all other points in signed distance field except the line of action. To sample contact points, one can simply draw samples from Eq. 2 and iterate from  $a$  to  $b$  until they reach a point that satisfies Eq. 11 and Eq. 12.

### B. Distribution on Surface Normals

Using Eq. 4 and Eq. 5, we can compute the mean of the gradient  $\mu_{\nabla}(x)$  and the covariance of the gradient  $\Sigma_{\nabla}(x)$  respectively. Thus we can compute the distribution around the surface normal for a given point in  $\mathcal{W}$ . We can now write

$$p(\mathbf{n}_i | \mathbf{c}_i = \gamma(t)) = p(\mathbf{n}_i | \mu(\gamma(t)), \Sigma(\gamma(t)))$$

One interesting effect of this technique is that we can now marginalize out the line of action model and visual what the surface normal distribution is along a given line of action. To our knowledge this is the first attempt to visual surface normals along a grasp plan. Marginalization can be performed as follows:

$$p(\mathbf{n}_i) = \int_a^b p(\mathbf{n}_i = \mathbf{v} | \mathbf{c}_i = \gamma(t)) p(\mathbf{c}_i = \gamma(t)) dt \quad (14)$$

Grasp metrics such as Ferrari-Canny require  $\mathbf{n}_i$  be normalized, or, equivalently, a member of the sphere  $\mathcal{S}^{d-1}$  [12]. To account for this we densely sample from the distribution  $p(\mathbf{n}_i)$  and project onto  $\mathcal{S}^{d-1}$ . In Fig.??, we visualize the theoretical distribution on  $\mathbf{n}_i$  calculated for a given GPIS and approach line of action.

### C. Expected Center of Mass

We recall the quantity  $P(sd(x) < 0) = \int_{-\infty}^0 p(sd(x) = s | \mu(x), \Sigma(x)) ds$  is equal to the probability that  $x$  is interior to the surface under the current observations. We assume that the object has uniform mass density and then  $P(sd(x) < 0)$  is the expected mass density at  $x$ . Then we can find the expected center of mass as:

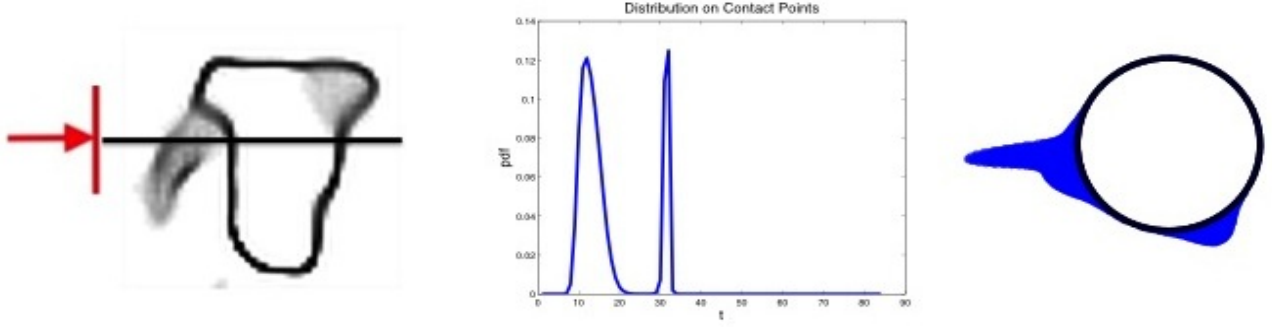


Fig. 5: (Left to Right): Line of action for a given gripper on an uncertain surface representing a measuring cup. Distribution  $p(c)$  as a function of  $t$ , the position along the line of action  $\gamma(t)$ . The two modes correspond to the different potential contact points, either the handle or the base of the cup. Lastly, the distribution on the surface normals (inward pointing) along  $\gamma(t)$  described by equation ??.

$$\bar{z} = \frac{\int_{\mathcal{W}} x P(sd(x) < 0) dx}{\int_{\mathcal{W}} P(sd(x) < 0) dx} \quad (15)$$

which can be approximated by sampling  $\mathcal{W}$  in a grid and approximating the spatial integral by a sum. Since this operation involves the entire SDF, one would want to use a low resolution grid for computational efficiency. We show the computed density and calculated expected center of mass for a marker in Fig. 6.

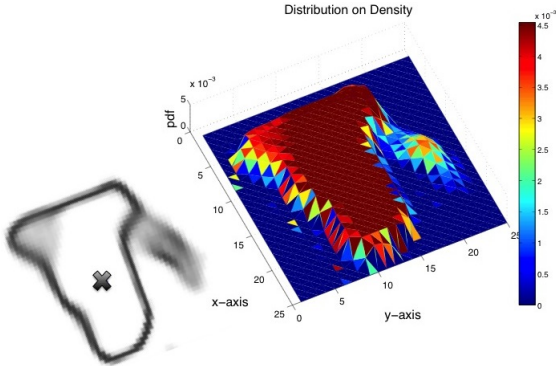


Fig. 6: Left: A surface with GPIS construction and expected center of mass (black X) Right: The distribution on the density of each point assuming uniform density

#### D. Complexity Analysis on Sampling From Grasps Distributions

After formally deriving each distribution, we can now sample along each line of action  $\gamma_i(t)$  from the joint distribution  $p(\mathbf{c}_i, \mathbf{n}_i | \gamma_i(t))$ , due to our independence assumption Eq. 9. This can be done by drawing samples from Eq. ?? and using our projection technique for the normal distribution.

Having a distribution along a line of action model allows us to sample from those instead of the joint distribution  $p(sd(\mathcal{W}))$ . Assuming the line of action is on the order of  $n$ , sampling from this distribution for a single grasp is  $O(n^3)$ . However, each proposed grasp plan  $\Gamma$  requires the distribution to be computed, so if we have  $T = |G|$  then the complexity is  $O(Tn^3)$ . In practice, this should be much smaller than  $O(n^6)$ .

### VIII. EXPERIMENTS

For the experiments below we used common household objects shown in Fig. ?? . We manually created a 25 x 25 grid, by tracing a pointcloud of the object on a table taken with a Primesense Carmine depth sensor. To accompany the SDF, we created an occupancy map, which holds 1 if the point cloud was observed and 0 if it was not observed, and a measurement noise map, which holds the variance 0-mean noise added to the SDF values. The parameters of the GPIS were selected using maximum likelihood on a held-out set of validation shapes. Our visualization technique follows the approach of [?] and consisted of drawing many shape samples from the distribution and blurring accordingly to a histogram equalization scheme.

We did experiments for the case of two hard contacts in 2-D, however our methods are not limited to this implementation. We drew random lines of actions  $\gamma_1(t)$  and  $\gamma_2(t)$  by sampling around a circle with radius  $\sqrt{2}n$  and sampling the circles origin, then projecting onto the largest inscribing circle in the workspace.

#### A. Multi-Arm Bandit Experiments

We consider the problem of selecting the best grasp plan,  $\Gamma^*$  out of a set  $G$ . For our experiments we look at selecting the best grasp out of a size of  $|G| = 1000$ . In Fig. 7, we plotted the simple regret for three of the shapes in our data set averaged over 100 runs and compare it to the Monte-Carlo method that randomly chooses a grasp plan to draw a sample from. We initialize both the Monte-Carlo and bandit technique by sampling each grasp 2 times this is to achieve the standard deviation. We draw samples from our calculated distributions  $p(g)$ . The most interesting thing is that the regret is minimized at least an order of magnitude faster than the naive approach for the shapes we considered, thus motivating the use of including observations as you select samples.

#### B. Sampling from Grasps Vs. Shape

We first tested 1000 grasp plans and sampled each one 5000 times and measured the RMS error between converged expected grasp plan qualities for sampling shape Eq. 6 vs. grasps Eq. 10 was 0.004. After confirming the distributions

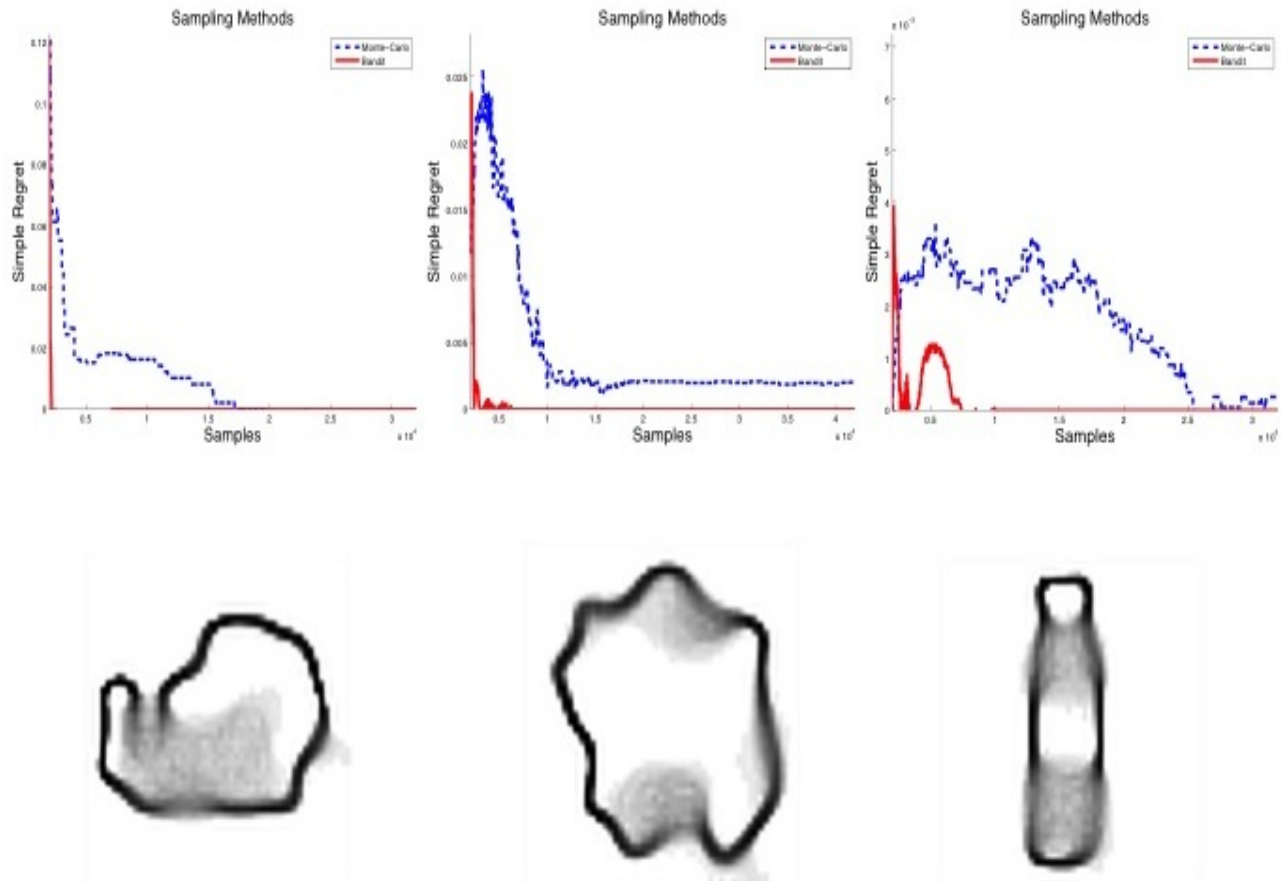


Fig. 7: Top: Comparison of two sampling methods. Red is the multi-arm bandit method that actively chooses which one to sample and Blue is the naive Monte-Carlo method. We measure their ability to converge in terms of simple regret averaged over 100 runs. We had them determine the best grasp when  $|G| = 1000$ . In all cases the bandit method converges at least a magnitude faster than the pure Monte-Carlo integration. Bottom: Three objects from our data set that we tested on (Tape, Loofa, Water Bottle), we used a visualization technique described in [?]

converged close to the same value, we show the computational complexity in Fig. 8 of the two methods for evaluating 100 grasps on an  $n \times n$  grid.

## IX. LIMITATIONS

Our budgeted multi-arm bandit approach appears promising, but we still do not know how well it will perform on 3D shapes and large scale grids. Future work will be building an efficient construction of GPIS to scale to 3D and test the bandit method there. While we have empirically seen the bandit method always converge to the correct value, our method only has a statistical guarantee of doing so. A low quality grasp plan could be found, albeit with a low probability.

Sampling from our distribution  $p(g)$  over  $p(S)$  yields a reduction in computational complexity, but only if the number of grasps one wants to evaluate remains small relative to  $n^3$ , techniques to ensure this could be to find locally optimal potential grasps using optimization approaches [?].

An additional problem is that we only have an expected center of mass and not a distribution on the center of mass. This might prove to be too expensive to compute, however recent work by Panahi et al. showed a way to bound the

center of mass for convex parts. Extension of his work to implicit surfaces could be of possible interest [28].

## X. CONCLUSION

Assessing grasp quality under shape uncertainty is computationally expensive as it often requires repeated evaluations of the grasp metric over many random samples. In this work, we proposed a multi-arm bandit approach to efficiently identify high-quality grasps under shape uncertainty. A key insight from our work is that uniformly allocating samples to grasps is inefficient, and we found that the Successive Elimination multi-arm bandit approach prioritizes evaluation of high-quality grasps while quickly pruning-out obviously poor grasps. A pre-requisite for applying a bandit approach is to formulate an efficient representation of how shape uncertainty affects grasp parameters and thus grasp quality. We modeled uncertainty with Gaussian Process Implicit Surfaces (GPIS) and derived the distribution of grasp parameters when a nominal grasp is applied to the GPIS. As a result, we were able to more efficiently sample from a distribution of grasps executions rather than the shape; leading to a complexity improvement of  $n^3$  in the resolution of the representation. We evaluated this theoretical model on a dataset of common objects and confirmed that: (1) the bandits approach always

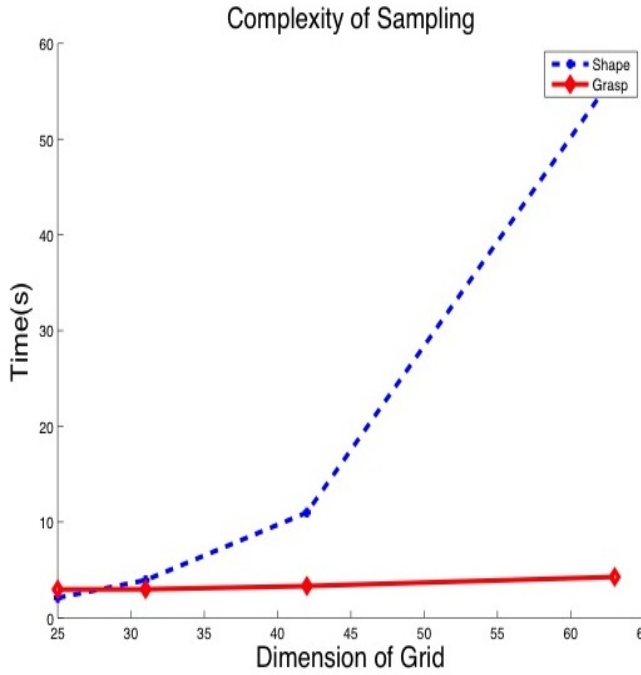


Fig. 8: Time it took to sample from 100 grasp distributions for a given resolution of the workspace. Blue line is sampling from  $p(sd(\mathcal{R}))$  or shapes and Red is sampling from  $p(g)$  or the calculated distribution on grasps. As you can see sampling from the calculated distributions scales much better.

converged to the best grasp in the candidate set, (2) it performs on average a magnitude faster than a naive uniform sampling approach in our experiments, and (3) sampling from the grasp contacts is 12x faster than sampling from shapes for a 64x64 grid.

## XI. FUTURE WORK

Our results are promising and they suggest many avenues of future work. First, we are actively working to apply this approach to 3D environments. In particular, we are interested in investigating the entire end-to-end pipeline of perception, grasp planning, and grasp execution. Next, in principle, our method can be applied to other representations of shape uncertainty such as perturbations on polygonal vertices [17] or splines [10]. It can further be applied to other grasp quality metrics. Next, we are investigating the application of recent results from reinforcement learning or bandit theory for analytic stopping criteria for the sampling. One particular challenge is that shape uncertainty is rarely parametric leading to problems with many recent bounds and optimality conditions.

## REFERENCES

- [1] R. Agrawal, "Sample mean based index policies with  $o(\log n)$  regret for the multi-armed bandit problem," *Advances in Applied Probability*, pp. 1054–1078, 1995.
- [2] S. Agrawal and N. Goyal, "Analysis of thompson sampling for the multi-armed bandit problem," *arXiv preprint arXiv:1111.1797*, 2011.
- [3] C. Andrieu, N. De Freitas, A. Doucet, and M. I. Jordan, "An introduction to mcmc for machine learning," *Machine learning*, vol. 50, no. 1-2, pp. 5–43, 2003.
- [4] P. Auer, N. Cesa-Bianchi, and P. Fischer, "Finite-time analysis of the multiarmed bandit problem," *Machine learning*, vol. 47, no. 2-3, pp. 235–256, 2002.

- [5] A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 1998.
- [6] D. Bergemann and J. Välimäki, "Bandit problems," Cowles Foundation for Research in Economics, Yale University, Tech. Rep., 2006.
- [7] S. Bubeck, R. Munos, and G. Stoltz, "Pure exploration in multi-armed bandits problems," in *Algorithmic Learning Theory*. Springer, 2009, pp. 23–37.
- [8] R. E. Caflisch, "Monte carlo and quasi-monte carlo methods," *Acta numerica*, vol. 7, pp. 1–49, 1998.
- [9] O. Cappé, A. Garivier, O.-A. Maillard, R. Munos, G. Stoltz, *et al.*, "Kullback-leibler upper confidence bounds for optimal sequential allocation," *The Annals of Statistics*, vol. 41, no. 3, pp. 1516–1541, 2013.
- [10] V. N. Christopoulos and P. Schrater, "Handling shape and contact location uncertainty in grasping two-dimensional planar objects," in *Intelligent Robots and Systems, 2007. IROS 2007. IEEE/RSJ International Conference on*. IEEE, 2007, pp. 1557–1563.
- [11] S. Dragiev, M. Toussaint, and M. Gienger, "Gaussian process implicit surfaces for shape estimation and grasping," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2011, pp. 2845–2850.
- [12] C. Ferrari and J. Canny, "Planning optimal grasps," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 1992, pp. 2290–2295.
- [13] K. Y. Goldberg and M. T. Mason, "Bayesian grasping," in *Robotics and Automation, 1990. Proceedings., 1990 IEEE International Conference on*. IEEE, 1990, pp. 1264–1269.
- [14] G. A. Hollinger, B. Englert, F. S. Hover, U. Mitra, and G. S. Sukhatme, "Active planning for underwater inspection and the benefit of adaptivity," *Int. J. Robotics Research (IJRR)*, vol. 32, no. 1, pp. 3–18, 2013.
- [15] K. Hsiao, M. Ciocarlie, and P. Brook, "Bayesian grasp planning," in *ICRA 2011 Workshop on Mobile Manipulation: Integrating Perception and Manipulation*, 2011.
- [16] E. Kaufmann, O. Cappé, and A. Garivier, "On bayesian upper confidence bounds for bandit problems," in *International Conference on Artificial Intelligence and Statistics*, 2012, pp. 592–600.
- [17] B. Kehoe, D. Berenson, and K. Goldberg, "Estimating part tolerance bounds based on adaptive cloud-based grasp planning with slip," in *Automation Science and Engineering (CASE), 2012 IEEE International Conference on*. IEEE, 2012, pp. 1106–1113.
- [18] —, "Toward cloud-based grasping with uncertainty in shape: Estimating lower bounds on achieving force closure with zero-slip push grasps," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 576–583.
- [19] J. Kim, K. Iwamoto, J. J. Kuffner, Y. Ota, and N. S. Pollard, "Physically-based grasp quality evaluation under uncertainty," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 3258–3263.
- [20] J. Laaksonen, E. Nikandrova, and V. Kyriki, "Probabilistic sensor-based grasping," in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*. IEEE, 2012, pp. 2019–2026.
- [21] T. L. Lai and H. Robbins, "Asymptotically efficient adaptive allocation rules," *Advances in applied mathematics*, vol. 6, no. 1, pp. 4–22, 1985.
- [22] B. León, S. Ulbrich, R. Diankov, G. Puche, M. Przybylski, A. Morales, T. Asfour, S. Moio, J. Bohg, J. Kuffner, and R. Dillmann, *Open-GRASP: A Toolkit for Robot Grasping Simulation*, ser. Lecture Notes in Computer Science. Springer Berlin / Heidelberg, 2010, vol. 6472, pp. 109–120.
- [23] Z. Li and S. S. Sastry, "Task-oriented optimal grasping by multifingered robot hands," *Robotics and Automation, IEEE Journal of*, vol. 4, no. 1, pp. 32–44, 1988.
- [24] O. Madani, D. J. Lizotte, and R. Greiner, "The budgeted multi-armed bandit problem," in *Learning Theory*. Springer, 2004, pp. 643–645.
- [25] J. Mahler, S. Patil, B. Kehoe, J. van den Berg, M. Ciocarlie, P. Abbeel, and K. Goldberg, "Gp-gpis-opt: Grasp planning under shape uncertainty using gaussian process implicit surfaces and sequential convex programming," in *Robotics and Automation (ICRA), 2015 IEEE International Conference on*. IEEE, 2015.
- [26] A. T. Miller and P. K. Allen, "Graspt! a versatile simulator for robotic grasping," *Robotics & Automation Magazine, IEEE*, vol. 11, no. 4, pp. 110–122, 2004.
- [27] D. F. Morrison, "Multivariate statistical methods. 3," *New York, NY. Mc*, 1990.
- [28] F. Panahi and A. F. van der Stappen, "Bounding the locus of the center of mass for a part with shape variation," *Computational Geometry*, vol. 47, no. 8, pp. 847–855, 2014.
- [29] K. B. Petersen, "The matrix cookbook."



- [30] C. Rasmussen and C. Williams, *Gaussian processes for machine learning*. MIT Press, 2006.
- [31] C. E. Rasmussen and H. Nickisch, "Gaussian processes for machine learning (gpml) toolbox," *The Journal of Machine Learning Research*, vol. 9999, pp. 3011–3015, 2010.
- [32] H. Robbins, "Some aspects of the sequential design of experiments," *Bulletin of the American Mathematical Society*, vol. 58, pp. 527–535, 1952.
- [33] M. Rothschild, "A two-armed bandit theory of market pricing," *Journal of Economic Theory*, vol. 9, no. 2, pp. 185–202, 1974.
- [34] R. Simon, "Optimal two-stage designs for phase ii clinical trials," *Controlled clinical trials*, vol. 10, no. 1, pp. 1–10, 1989.
- [35] E. Solak, R. Murray-Smith, W. E. Leithead, D. J. Leith, and C. E. Rasmussen, "Derivative observations in gaussian process models of dynamic systems," 2003.
- [36] D. L. St-Pierre, Q. Louveaux, and O. Teytaud, "Online sparse bandit for card games," in *Advances in Computer Games*. Springer, 2012, pp. 295–305.
- [37] F. Stulp, E. Theodorou, J. Buchli, and S. Schaal, "Learning to grasp under uncertainty," in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE, 2011, pp. 5703–5708.
- [38] N. Vahrenkamp, T. Asfour, and R. Dillmann, "Simo: A simulation and motion planning toolbo for c+."
- [39] R. Weber *et al.*, "On the gittins index for multiarmed bandits," *The Annals of Applied Probability*, vol. 2, no. 4, pp. 1024–1033, 1992.
- [40] J. Weisz and P. K. Allen, "Pose error robust grasping from contact wrench space metrics," in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 557–562.
- [41] O. Williams and A. Fitzgibbon, "Gaussian process implicit surfaces," *Gaussian Proc. in Practice*, 2007.
- [42] Y. Zheng and W.-H. Qian, "Coping with the grasping uncertainties in force-closure analysis," *Int. J. Robotics Research (IJRR)*, vol. 24, no. 4, pp. 311–327, 2005.