

Dex-Net 1.0: A Cloud-Based Network of 3D Object Models Based on Deep Learning for Object Similarity and Multi-Armed Bandits for Robust Grasp Planning

v4 [2015-09-02 17:35]

Jeffrey Mahler¹, Florian Pokorny¹, Brian Hou¹, Mel Roderick¹, Michael Laskey¹, Kai Kohlhoff³,
Torsten Kroeger³, James Kuffner³, Ken Goldberg²

[TODO: POSSIBLY WILL INCLUDE MATTHIEU AUBRY,
ALYOSHA EFROS, AND PIETER ABBEEL AS COAUTHORS]
[TODO: UPDATE NUMBER OF MODELS IN TITLE: I WANT
IT TO REFLECT EXACTLY HOW MANY WE USE IN EXPERI-
MENTS]

Abstract— Advances in Big Data and distributed computing, combined with very large datasets of images and human speech, have produced a series of data-driven results that surpass decades of research in analytic methods and, in some cases, that surpass human ability. For example, deep neural networks can now learn to play video games, recognize faces, and translate between languages. This raises the question: can Big Data and computation produce analogous advances in robot grasping and manipulation? We introduce the Dexterity-Network (Dex-Net) 1.0, a Cloud-based dataset of 3D object models labelled with parallel-jaw grasps and similarity between objects and grasps based on multi-view Convolutional Neural Networks (CNNs) and grasp heightmaps. We use Dex-Net in Multi-Armed Bandit (MAB) algorithms to quickly find a grasp with high probability of force closure on an object from a set of 250 candidates. We extend previous MAB models to predict grasp quality using similarity to prior objects and grasps in Dex-Net using Continuous Correlated Beta Processes (CCBPs). Our initial results suggest that using prior objects from Dex-Net can accelerate grasp planning by up to $10\times$. Code, data, and additional information can be found at: [TODO: URL]

I. INTRODUCTION

Consider a robot fulfilling orders in a warehouse. The robot must quickly plan grasps for consumer products to fulfill orders. Furthermore, the robot may not precisely know the state of the environment, such as the pose or frictional properties of object, due to sensor imprecision and missing data (e.g. occlusions in point clouds). This poses a challenge for many existing grasp planning methods that are based on precise knowledge of contact locations and surface normals [12], [19] due to the distribution on possible environment configurations.

In vision and speech, advances in Big Data and distributed computing, combined with datasets of millions of examples



Fig. 1: Sample of 3D mesh models from the Dex-Net dataset. The dataset currently includes over 10,000 models from laser-scanned datasets such as the KIT object database [36] and the Yale-CMU-Berkeley object set [7], and synthetic datasets such as 3DNet[73], ModelNet [74], and the SHREC 2014 object retrieval challenge dataset [49].

such as ImageNet [13] and the Fisher corpus [11], have produced impressive results [26], [27], [41] that surpass those obtained from decades of research on analytic methods. This raises the question: will machine learning of grasps for vast numbers of possible object poses, object shapes, environment configurations, etc., exhibit scaling effects similar to those observed in computer vision and speech recognition?

We introduce Dex-Net 1.0, a dataset of 3D models and grasps with similarity metrics, and a system for actively acquiring statistical models of grasp quality across the network. The Dex-Net dataset contains representative 3D models of objects that could be encountered in warehousing or the home. Dex-Net contains laser-scanned 3D mesh models from the KIT object database [36], the Amazon Picking Challenge objects, BigBIRD [66], and YCB [7] to reflect physical objects commonly used for benchmarking in grasping research. The dataset also includes synthetic 3D mesh models from shape classification datasets such as 3DNet [73], ModelNet [74], and the SHREC 2014 large scale object retrieval challenge [49] to scale the number of models. Fig. 1 shows a sample of the objects in the dataset. We measure object similarity based the distance between feature vectors generated from multi-view Convolutional Neural Networks (CNNs) [70], a state-of-the art method in shape retrieval. We measure similarity between grasps by the pose of the gripper and local surface “heightmaps” extracted at location of contact with an object [28], [35]. Our system uses Google Cloud Platform to distribute grasp quality evaluations across hundreds of virtual machines and to store millions of grasps across the objects in Dex-Net.

In this work, we use Dex-Net to study the convergence rate of Multi-Armed Bandit (MAB) algorithms for finding a grasp with high probability of force closure under object pose, gripper pose, and friction coefficient uncertainty from a set of 250 candidate grasps on objects from Dex-Net. To do so, we extend the MAB framework of Laskey et al. [44] to model correlations between the probability of force closure for similar grasps on similar 3D objects using Continuous Correlated Beta Processes (CCBPs) [21], [55]. CCBPs predict a belief distribution for the quality of each grasp based similarity to data in Dex-Net that can be efficiently updated after observing the quality of a grasp for a sampled object pose, gripper pose, and friction coefficient. Our results suggest that CCBPs bootstrapped with data from Dex-Net can accelerate grasp planning by up to 10 \times .

II. RELATED WORK

For a survey of the substantial literature on grasping, see Prattichizzo and Trinkle [62]. Research on grasp planning has focused on finding grasps by maximizing an analytic grasp quality metric based on wrench space (WS), such as the ability to resist external perturbations to the object based on the grasp WS [19], [54], object WS [52], [61], or task WS [43], [51]. WS metrics have been used to synthesize grasps for known objects using sampling-based optimization in software tools such as GraspIt! [54] or OpenGRASP [47]. However, analytic WS metrics have been criticized [1], [72] for not being robust to variations in grasp acquisition, such as object shape, pose, material properties, and locations of contact [8], [17], [72], [75] and for not taking into account the dynamics of the the grasp [62].

The need to recompute WS metrics to select grasps for every new object motivated research on grasp synthesis by transferring grasps from a database of exemplar objects. Li and Pollard [50] generated grasps by matching object shapes to human hand postured in a database. Goldfelder et al. [24] developed the Columbia grasp database, a database of 1,814 distinct models and over 200,000 force closure grasps generated using the Eigengrasp planner in GraspIt! [12]. The authors later used synthetic partial depth maps of objects in the database to match robot sensor data to precomputed grasps, using the Iterated Closest Point (ICP) algorithm to align the coordinate frames of the depth maps [25], [23]. Kehoe et al. [39] created a Cloud-based system to transfer grasps evaluated by probability of force closure on objects in a database to a physical robot by indexing the objects with the Google Goggles object recognition engine. Recent research has also studied grasp transfer from objects of the same category by warping contacts between corresponding points on a shape surface and using local rigid alignment and contact interpolation [30], [68] or by interpolating grasps and shapes over a vector space representation called a Grasp Moduli Space [58], [59].

Another line of research focused on making analytic grasp metrics robust to imprecision in perception and control [22], [69], [75]. Brook, Ciocarlie, and Hsiao [6], [31] developed a Bayesian framework to evaluate both the expected epsilon

quality and the probability of physical success on a PR2 given uncertainty in object identity, object pose, and gripper positioning on deterministic mesh and point cloud models. Weisz et al. [72] found that grasps ranked by probability of force closure subject to perturbations of object pose in simulation were empirically more successful on a physical robot than grasps planned using deterministic WS metrics. Kim et al. [40] planned grasps using the expected epsilon quality metric [19] under dynamics and uncertainty in pose, and found that the robust metric has a higher correlation with physical grasp success. Recent research has also studied grasping under object shape uncertainty resulting from imprecision of object segmentations in images [9], part tolerancing in manufacturing [37], [38], [57], or missing and noisy data from depth sensors such as the Kinect modeled with Gaussian process implicit surfaces [18], [53], [44].

Recent research has also studied synthesizing grasps by ranking grasps according statistical models learned from human annotations or physical execution [3]. Saxena et al. [33], [65] used a logistic regression classified to predict grasp affordances in images from human annotated training data. Lenz et al [46] used deep learning to detect bounding boxes for parallel-jaw grasps in color and depth images, which was extended to real time by Redmon and Angelova [63]. Herzog et al. [28], [29] extracted "heightmaps" of local object curvature from human demonstrated grasps, construct a library of heightmap templates, and match new sensor data to templates to select grasps similar to the demonstrations. Detry et al. [15] created a low-dimensional representation of object parts and cluster object parts that are grasped similarly to form a shape library of prototypical grasp parts, and show that this representation can be transferred to real sensor data [14]. Kappler et al. [35] trained a deep neural network to predict grasp success for a Barrett hand measured by human annotations and the results of simulations on a database of synthetic pointclouds of objects. Deep learning [41] has also been used in robotics for learning visuomotor policies for specific manipulation tasks [48] and recurrent control policies for cutting fruits vegetables from joint angles and end-effector forces [45]. In comparison, we learn a model to predict a Bayesian distribution on the probability of force closure for a grasp on an object based on similarity to a set of prior grasps and objects in a database, and use our model to actively decide the next grasp to evaluate using Multi-Armed Bandits.

Our work is also closely related to research on actively selecting grasps for building a statistical model of grasp quality from fewer examples. Several works have searched for successful grasps using belief space planning to minimize uncertainty [32], [34], [20], but without use of an explicit grasp quality metric to guide the search. Kehoe et al. [37] proposed iterative pruning, an algorithm for evaluating the probability of force closure for a set of candidate grasps while discarding grasps known to have poor quality. Detry et al. [16] estimated a full continuous density function over grasp poses using kernel density estimation and adaptively acquired samples by pruning unsuccessful grasps. Kroemer

et al. [42] developed a reinforcement learning approach to grasp selection based on seeding hypotheses via imitation learning and Gaussian process upper-confidence bounds for active grasp acquisition. Boularias et al. [4], [5] used a Gaussian process Bayesian Optimization model for selecting grasps on cluttered piles of rocks. Salaganicoff et al. [64] used active learning to decide the next grasp to execute on a physical robot while learning a predictive model of empirical success from range sensors. Similarly, Montesano and Lopes [55] used Continuous Correlated Beta Processes [21] to actively acquire grasp executions on a physical robot, measuring correlations from the responses to a bank of image filters designed to detect grasp affordances such as edges. Recently, Laskey et al. [44] showed that Multi-Armed Bandit (MAB) algorithms can be used to accelerate the identification of grasps with high probability of force closure under uncertainty in shape, pose, and friction in 2D. MAB algorithms trade off gaining information about grasps that have been sampled fewer times with exploiting the grasp with the highest estimated quality given the past samples. In this work we extend the model of Laskey et al. [44] to 3D and to utilize similarity between grasps across objects from Dex-Net with CCBPs to further reduce the number of samples needed to converge to a grasp with high quality.

III. DEFINITIONS AND PROBLEM STATEMENT

Given a previously unknown object, we consider the problem of finding the parallel-jaw grasp with the maximum expected grasp quality according to a binary grasp quality metric such as force closure under uncertainty in object pose, gripper pose, and friction coefficient.

A. Grasp and Object Model

Our grasping model is illustrated in Fig. 2. Let $\mathbf{g} = [\mathbf{x}, \mathbf{w}]^T$ be a parallel-jaw grasp parameterized by the centroid of the jaws in 3D space $\mathbf{x} \in \mathbb{R}^3$ and an approach direction, or axis, $\mathbf{v} \in \mathbb{S}^2$. We assume that the jaws are opened to their maximal width $w \in \mathbb{R}$ before closing on the object. Let $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ be a signed distance function (SDF) representing an object's geometry [53], [56], which is zero on the object surface, positive outside the object, and zero inside the object. We assume that the object is specified in units of meters and the object center of mass $\mathbf{z} \in \mathbb{R}^3$ is known. Let $\mathcal{G} = \{[\mathbf{x}, \mathbf{w}]^T | \mathbf{x} \in \mathbb{R}^3, \mathbf{v} \in \mathbb{S}^2\}$ denote the space of all grasps and $\mathcal{H} = \{\mathcal{O} = \{\mathbf{z}, f(\cdot)\} | \mathbf{z} \in \mathbb{R}^3, f \in \mathcal{F}\}$ denote the space of all objects, where \mathcal{F} is the space of all SDFs. Our joint space of grasps and objects, or Grasp Moduli Space [59], is $\mathcal{M} = \mathcal{G} \times \mathcal{H}$.

B. Sources of Uncertainty

We assume known uncertainty in object pose, gripper pose, and friction coefficient resulting from sensing uncertainties and missing data following the model of Laskey et al. [44]. Let ξ be Gaussian uncertainty in object pose with mean $\bar{T} \in SE(3)$ and covariance Σ_ξ following the model of Barfoot and Furgale [2]. Let ν be uncertainty in gripper

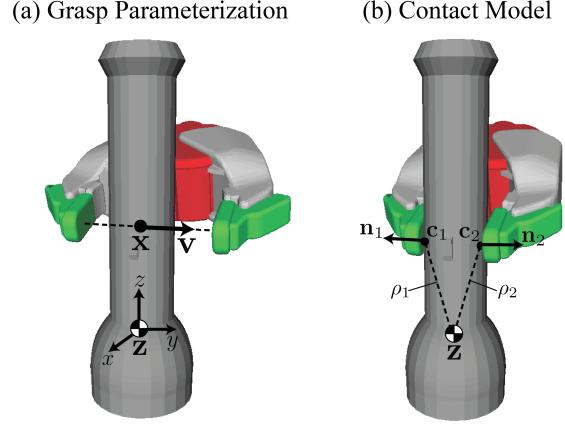


Fig. 2: (Left) We parameterize parallel-jaw grasps by the centroid of the jaws $\mathbf{x} \in \mathbb{R}^3$ and approach direction, or direction along which the jaws close, $\mathbf{v} \in \mathbb{S}^2$. The grasp parameters \mathbf{x} and \mathbf{v} are specified with respect to a coordinate frame located at the object center of mass \mathbf{z} and oriented along the principal directions of the object. (Right) The jaws are closed until contacting the object surface at locations $\mathbf{c}_1, \mathbf{c}_2 \in \mathbb{R}^3$, at which the surface has normals $\mathbf{n}_1, \mathbf{n}_2 \in \mathbb{S}^2$. The contacts are used to compute the moment arms $\rho_1 = \mathbf{c}_1 - \mathbf{z}$ and $\rho_2 = \mathbf{c}_2 - \mathbf{z}$. From these parameters we can derive the theoretical forces and torques that the gripper can apply to the object.

pose distributed as a Gaussian with mean $\mu_\nu \in \mathcal{G}$ and covariance Σ_ν . Let γ denote uncertainty in friction coefficient distributed as a Gaussian with mean $\mu_\gamma \in \mathbb{R}$ and covariance Σ_γ . These sources of uncertainty might be estimated based on prior models of repeatability in gripper actuation or from the posterior variance in object pose or friction coefficient in estimates based on sensor data.

C. Contact Model

Our contact model is illustrated in the right panel of Fig. 2. Given a grasp \mathbf{g} with width w on an object \mathcal{O} and samples of object pose, gripper pose, and friction, let $\mathbf{c}_i \in \mathbb{R}^3$ for $i \in 1, 2$ denote the 3D location of contact between the i -th gripper jaw and surface. We can compute the contacts using

$$\begin{aligned} \mathbf{c}_1 &= \min_{t \geq 0} t \text{ such that } |f(\mathbf{x} + (t - w/2)\mathbf{v})| < \epsilon \\ \mathbf{c}_2 &= \min_{t \geq 0} t \text{ such that } |f(\mathbf{x} - (t - w/2)\mathbf{v})| < \epsilon \end{aligned}$$

where $\epsilon > 0$ is a user-specified surface resolution [53].

Given contact points, let $\mathbf{n}_i = \nabla f(\mathbf{c}_i)/\|\nabla f(\mathbf{c}_i)\|_2$ denote the surface normal at contact \mathbf{c}_i with tangent vectors $\mathbf{t}_{i,1}, \mathbf{t}_{i,2} \in \mathbb{S}^2$. To compute the forces that each contact can apply to the object for a given friction coefficient γ , we discretize the friction cone at contact i into a discrete set of l facets with vertices $\mathcal{F}_i = \{\mathbf{f}_{i,j} = \mathbf{n}_i + \gamma \cos(\frac{2\pi j}{l}) \mathbf{t}_{i,1} + \gamma \sin(\frac{2\pi j}{l}) \mathbf{t}_{i,2} \text{ for } j = 1, \dots, l\}$ [60]. Each force $\mathbf{f}_{i,j}$ can exert a corresponding torque $\tau_{i,j} = \mathbf{f}_{i,j} \times \rho_i$ where $\rho_i = (\mathbf{c}_i - \mathbf{z})$ is the moment arm at contact i . To enable force closure with two contacts in 3D we assume a soft contact model, under which each contact \mathbf{c}_i exerts an additional wrench $\mathbf{w}_{i,l+1} = [\mathbf{0}, \mathbf{n}_i]^T$ [62]. Thus the set of all wrenches that can be applied by a grasp \mathbf{g} under our model is $\mathcal{W} = \{\mathbf{w}_{i,j} = [\mathbf{f}_{i,j}, \tau_{i,j}]^T | i = 1, 2 \text{ and } j = 1, \dots, l + 1\}$.

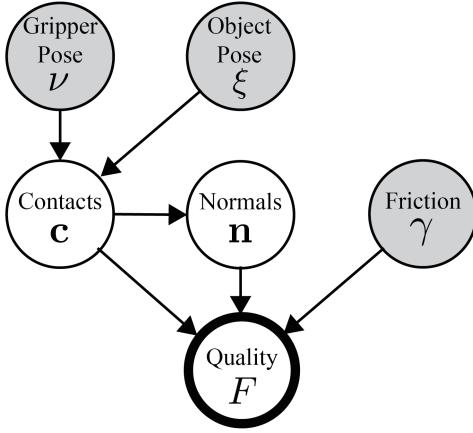


Fig. 3: A graphical model describing the relationships between known uncertain parameters (shaded gray) such as object pose and the force closure random variable F .

D. Quality Metric

In this work we use force closure, or the ability to resist external force and torques in arbitrary directions [19], as our grasp success metric. We acknowledge that force closure assumes the ability to actuate the gripper to apply arbitrary forces in the contact friction cone [], and that some evidence suggests that deterministic force closure is not always a strong predictor of physical grasp success [1]. However, we use the probability of force closure because of it has shown promise in physical experiments [40], [72], and it is relatively inexpensive to evaluate when compared to human labels or physical executions, allowing us to better study the effects of large amounts of data.

Let $F \in \{0, 1\}$ denote the occurrence of force closure. To compute force closure for a grasp $\mathbf{g} \in \mathcal{G}$ on object $\mathcal{O} \in \mathcal{H}$ given samples of object pose ξ , gripper pose ν , and friction coefficient γ , we first compute the set of possible contact wrenches \mathcal{W} . Then $F = 1$ if the origin lies within the convex hull of the contact wrench set, $\mathbf{0} \in \text{Conv}(\mathcal{W})$ [72]. A graphical model describing the relationship between these quantities is given in Fig. 3.

E. Objective

The probability of force closure for a grasp \mathbf{g} on object \mathcal{O} is

$$P_F(\mathbf{g}, \mathcal{O}) = \mathbb{P}(F = 1 \mid \mathbf{g}, \mathcal{O}, \xi, \nu, \gamma).$$

We are interested in finding the candidate grasp that maximizes the probability of force closure $P_F(\mathbf{g})$ [40], [44], [53], [72] subject to these sources of uncertainty over a budgeted maximum number of samples T . To perform this as quickly as possible we formulate this as a maximization over the sum of the true P_F for the grasps sampled at iteration t , where $I(i)$ denotes the grasp selected at time t [44]:

$$\underset{\mathbf{g}_{I(1)}, \dots, \mathbf{g}_{I(T)} \in \mathcal{G}}{\text{maximize}} \sum_{t=1}^T P_F(g_{I(t)}). \quad (\text{III.1})$$

As the maximization over the continuous space \mathcal{G} is computationally expensive, past work has solved this objective using a discrete set of K candidate grasps $\Gamma =$

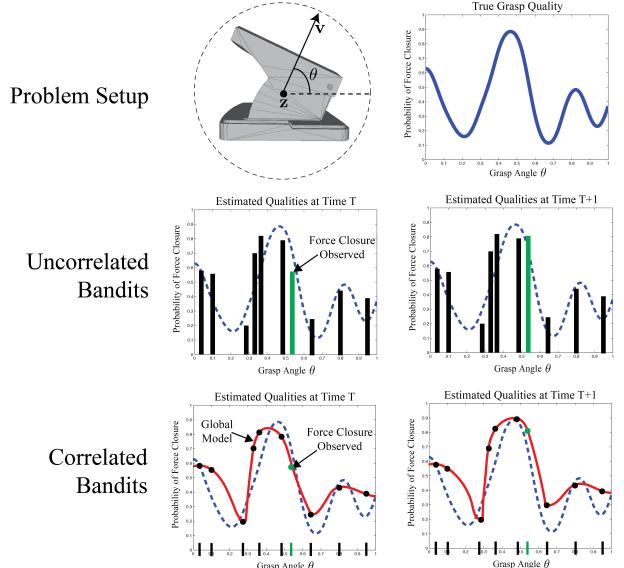


Fig. 4: (Top Left) Consider a set of grasps with fixed center at an object center of mass and approach direction sampled along a one-dimensional circle. (Top Right) The true probability of force closure P_F as a function of the angle of the grasp axis. (Middle Left) Uncorrelated MAB models maintain an estimate of P_F independently for each grasp, indicated by bars and their height. On iteration T of a MAB algorithm, the evaluated grasp (green bar) is force closure for samples . (Middle Right) Only the estimate of P_F for the sampled grasp is updated on iteration $T + 1$, and the grasp with highest estimated P_F remains suboptimal. (Bottom Left) Correlated MAB models maintain a global predictive function (red) of P_F for any possible grasp but may select grasps from a set of discrete candidates (indicated by dots). (Bottom Right) On iteration $T + 1$, the evaluated grasp (green dot) is sampled returns force closure and the global model is updated, increasing the estimated P_F for “nearby” grasps. Although a suboptimal grasp was sampled, the global model now correctly predicts the optimal grasp in the set.

$\{\mathbf{g}_1, \dots, \mathbf{g}_K\}$ where Γ is set by sampling grasp centers from a Gaussian with mean at the object center [44] or using heuristics specific to parallel-jaw grippers such as antipodal [53]. However, even with a discrete set of candidates, fully evaluating $P_F(\mathbf{g})$ for any grasp requires an expensive integration over possible contact locations between the grasp and surface and the surface normals at these contacts [53]. Thus, past research has evaluated the probability of force closure using Monte-Carlo integration [37], [38], [72], approximating the expression by minimizing uncertainty at the contact locations [53], and using Multi-Armed Bandit (MAB) algorithms to jointly evaluate $P_F(\mathbf{g})$ using Monte-Carlo integration while allocating samples to more promising grasps. In this work, we extend the MAB model of [44] to model similarities between grasps and prior objects in Dex-Net to further accelerate convergence.

IV. CORRELATED MULTI-ARMED BANDIT MODEL

Finding an approximate solution of Equation III.1 with Multi-Armed Bandits requires a predictive model of the probability of force closure for each grasp. We use Continuous Correlated Beta Processes (CCBPs) to model force closure for a grasp as a Bernoulli random variable with correlations across grasps and objects [21], [55]. Fig. 4 illustrates how using a correlated model can lead to accelerated convergence of bandit algorithms.

A. Continuous Correlated Beta Processes

Continuous Correlated Beta Processes (CCBPs) were independently proposed by Goetschalckx et al. [21] and Montesano and Lopes [55] to model correlations between the Bernoulli random variables in a Beta-Bernoulli process, which may lead to faster convergence in MAB problems [10]. Such correlations may exist when the Bernoulli random variables depend on common latent factors. For example, two grasps may have similar P_F when they contact the same shape in similar locations or the surface geometry at the contact points is similar, such as a grasp on the handle of a mug and the handle of a teapot.

We are interested in solving Equation III.1 for a new object \mathcal{O}_i . Let \mathcal{G}_i be a set of K candidate grasps for \mathcal{O}_i , $\Gamma_i = \{\mathbf{g}_{i,k}\}_{k=1}^K$. We define $F_{j,\ell} = F(\mathbf{g}_{j,\ell}) \in \{0, 1\}$ an occurrence of force closure on an evaluation of any grasp $\mathbf{g}_{j,\ell}$ on object \mathcal{O}_j using samples of object pose, gripper pose, and friction, as described in Section III-D. Then let $N_{j,\ell}$ be the number of total evaluations of for grasp $\mathbf{g}_{j,\ell}$, and define $S_{j,\ell} = \sum_{a=1}^{N_{j,\ell}} F_{j,\ell}$ to be the total number of times force closure has been observed for any grasp. Furthermore let $\mathcal{D} = \{N_{j,\ell}, S_{j,\ell}, \mathcal{Y}_{j,\ell} = (\mathbf{g}_{j,\ell}, \mathcal{O}_j) | j = \{1, \dots, M\}, \ell \in \{1, \dots, K\}\}$ be a discrete set of M known objects stored in a database, each with a discrete set of K candidate grasps $\Gamma_j = \{\mathbf{g}_{j,\ell}\}_{\ell=1}^K$ and $N_{j,\ell}$ and $S_{j,\ell}$ for each object.

We model force closure for each candidate grasp, $F_{i,k}$, as a Bernoulli random variable with probability of success $\theta_{i,k} = P_F(\mathbf{g}_{i,k})$. Since we do not know the value of $\theta_{i,k}$, we maintain a belief distribution for each $\theta_{i,k}$ based on our prior belief about the likelihood of force closure. In CCBPs, the belief distribution on the Bernoulli parameter $\theta_{i,k}$ is the Beta distribution, which is specified by shape parameters $\alpha > 0$ and $\beta > 0$:

$$\text{Beta}(\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta_k^{\alpha-1} (1 - \theta_k)^{\beta-1}$$

A CCBP estimates the shape parameters for a grasp and object $\mathcal{Y}_{i,k} = (\mathbf{g}_{i,k}, \mathcal{O}_i) \in \mathcal{M}$ using a normalized kernel function $k(\mathcal{Y}_p, \mathcal{Y}_q) : \mathcal{M} \times \mathcal{M} \rightarrow [0, 1]$ that measures similarity between a pair of grasps and objects from the Grasp Moduli Space \mathcal{M} . The kernel approaches 1 as the arguments become increasingly similar and approaches 0 as the arguments become dissimilar. In this work we use a unit-bandwidth squared exponential kernel

$$k(\mathcal{Y}_p, \mathcal{Y}_q) = \exp \left(-\frac{1}{2} \sum_{m=1}^{N_f} w_m^2 \|\varphi_m(\mathcal{Y}_p) - \varphi_m(\mathcal{Y}_q)\|_2^2 \right)$$

where $\varphi_m : \mathcal{M} \rightarrow \mathbb{R}^{d_m}$ for $m = 1, \dots, N_f$ are feature mappings for a grasp and object to a d_m -dimensional Euclidean space and $w_m \in \mathbb{R}$ are weights scaling the relative contribution of the feature maps to the similarity metric.

Before the first iteration of the MAB algorithm, we update our belief for each candidate grasp to its similarity to all grasps and objects from the database \mathcal{D} as measured by the

kernel [21]:

$$\alpha_{i,k,0} = \alpha_0 + \sum_{j=1}^M \sum_{\ell=1}^K k(\mathcal{Y}_{i,k}, \mathcal{Y}_{j,\ell}) S_{j,\ell} \quad (\text{IV.1})$$

$$\beta_{i,k,0} = \beta_0 + \sum_{j=1}^M \sum_{\ell=1}^K k(\mathcal{Y}_{i,k}, \mathcal{Y}_{j,\ell})(N_{j,\ell} - S_{j,\ell}) \quad (\text{IV.2})$$

where α_0 and β_0 are prior parameters for the Beta distribution [44]. Upon observing $F_{i,k}$ for grasp $\mathbf{g}_{i,k}$ on iteration t , we update our belief $\theta_{j,k}$ for all other grasps on object \mathcal{O}_i by [21]:

$$\alpha_{j,k,t} = \alpha_{j,k,t-1} + k(\mathcal{Y}_{j,k}, \mathcal{Y}_{i,k}) F_{i,k} \quad (\text{IV.3})$$

$$\beta_{j,k,t} = \beta_{j,k,t-1} + k(\mathcal{Y}_{j,k}, \mathcal{Y}_{i,k})(1 - F_{i,k}) \quad (\text{IV.4})$$

Intuitively, this allows observations of one grasp to constitute fractional observations of similar grasps.

B. Feature Mappings

We use a set of feature mappings $\varphi_1, \dots, \varphi_{N_f}$ to capture similarity based on grasp parameters, local surface geometry, and global object shape.

1) *Grasp Parameters*: Since grasps may be correlated across small changes to the grasp center and approach direction on a single shape, we use the feature maps $\varphi_x(\mathcal{Y}) = \mathbf{x}$ and $\varphi_v(\mathcal{Y}) = \mathbf{v}$, the identity transformations for the grasp center and axis. Additionally we use the moment arm feature map $\varphi_\rho(\mathcal{Y}) = [\|\rho_1\|_2, \|\rho_2\|_2]^T \in \mathbb{R}^2$ because points on the object surface with larger moment arms will move greater distances under object orientation uncertainty. However, similarity between grasp centers and approach directions may not indicate similar P_F if the surfaces near the contact locations are dissimilar. For example, consider two grasps that contact a box near the corner. The two grasps may have similar parameters and moment arms but may contact the surface on different sides of the edge, resulting in different surface orientations at the contact points.

2) *Grasp Heightmaps*: Section ?? We also use a variant of the grasp heightmap features of Herzog et al. [28] and Kappler et al. [35]. Let $d_h \in \mathbb{Z}$ be the number of pixels along each dimension of the heightmap, let $\mathcal{P} = \{-d_h, \dots, d_h\}$ be the row / column pixel indices for the heightmap, let $\delta \in \mathbb{R}$ be the resolution of the image pixels in meters, and let $r \in \mathbb{R}$ be a minimum projection distance. Furthermore, let $\mathbf{c}_i, i = 1, 2$ be a contact point for grasp \mathbf{g} and let $\mathbf{t}_1, \mathbf{t}_2$ be two orthogonal unit vectors to the grasp approach direction \mathbf{v} . Our heightmap at contact i , $\mathbf{h}_i : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$, maps discrete locations along the tangent plane specified by $\mathbf{t}_1, \mathbf{t}_2$ to the distance to the surface along the grasp axis \mathbf{v} . To compute the heightmap value at pixel $u, v \in \mathcal{P}$, we first compute the 3D location of the pixel on the plane $\mathbf{p}_i(u, v) = \mathbf{c}_i + \delta u \mathbf{t}_1 + \delta v \mathbf{t}_2$. Then we assign the heightmap value

$$\mathbf{h}_i(u, v) = \min_{t \geq -r} t \text{ such that } |f(\mathbf{p}_i(u, v) + s_k \mathbf{v})| < \epsilon$$

where $s_1 = -1$, $s_2 = 1$, ϵ is the surface threshold used for finding contacts, and f is the SDF of object \mathcal{O} . Finally, we make the heightmap rotationally invariant by rotating it to

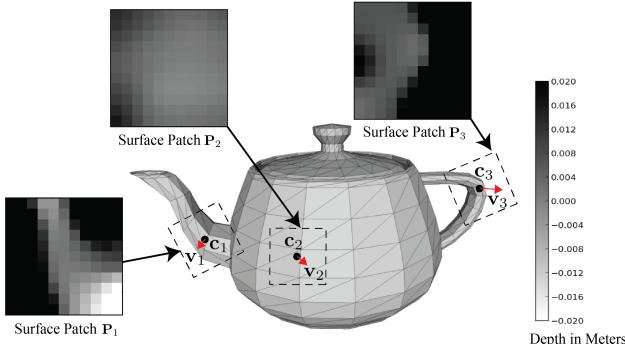


Fig. 5: Three local surface heightmaps extracted on a teapot. Each heightmap is “rendered” along the grasp axis at each contact point and oriented by the local directions of maximum variation in the heightmap.

align the axes with the eigenvectors of the weighted covariance matrix of contact points generating the heightmap as described in [71]. Our full feature vector for the heightmaps is $\varphi_h(\mathbf{g}_i, \mathcal{O}_i) = [\mathbf{h}_1, \mathbf{h}_2]^T$. [TODO: THIS IS SLIGHTLY MORE COMPLICATED BUT TAKES A BIT TO EXPLAIN; BASICALLY THE ORDERING OF CONTACTS DOESN’T MATTER HERE]

Fig. 5 illustrates local surface patches extracted by this procedure.

3) Global Features: While force closure is based on local surface properties of an object around contact points, measuring global similarity may be similar. It may be useful for two reasons. First, the sums in Equations IV.1 and IV.2 may be expensive when the database is large, and global object similarity can be used to sum over only the objects that are grasped. Second, the heightmaps may not capture all possible contact points for a grasp under perturbations in object and gripper pose.

We measure global similarity between objects using a feature mapping derived from multi-view Convolutional Neural Networks (CNNs) [70], a state-of-the-art method for 3D shape retrieval. Fig. 6 illustrates our method. Let R be the maximum dimension of the object bounding box. We first render every object on a white background in a total of N_c virtual camera views oriented toward the object center and discretized on a viewing sphere along angle increments $\delta_\theta = \frac{\pi}{N_c}$ and $\delta_\varphi = \frac{\pi}{2N_c}$ and radii $r = R, 2R$. Then we train a deep CNN with the architecture of AlexNet [41] to predict a class label for the rendered images based on known class labels for the 3D models using Stochastic Gradient Descent (SGD). Next, we pass each of the N_c views of each object through the finetuned CNN and max-pool the output of the fc7 layer. We finally use Principal Component Analysis (PCA) to reduce the max-pooled output to 400 dimensions to 100 dimensions. This yields a representation $\psi(\mathcal{O}) \in \mathbb{R}^{100}$ for each object, and thus our final feature map is $\varphi_g(\mathcal{Y}) = \psi(\mathcal{O})$. In our implementation, we trained the multi-view CNN on rendered images of 171 classes of 3D models from the SHREC 2014 dataset [49] for 500,000 iterations of SGD, which had a test accuracy of 76%.

C. Optimizing Feature Weights

One remaining issue is the selection of the relative feature weights w_1, \dots, w_{N_f} . We select the weights that minimize the

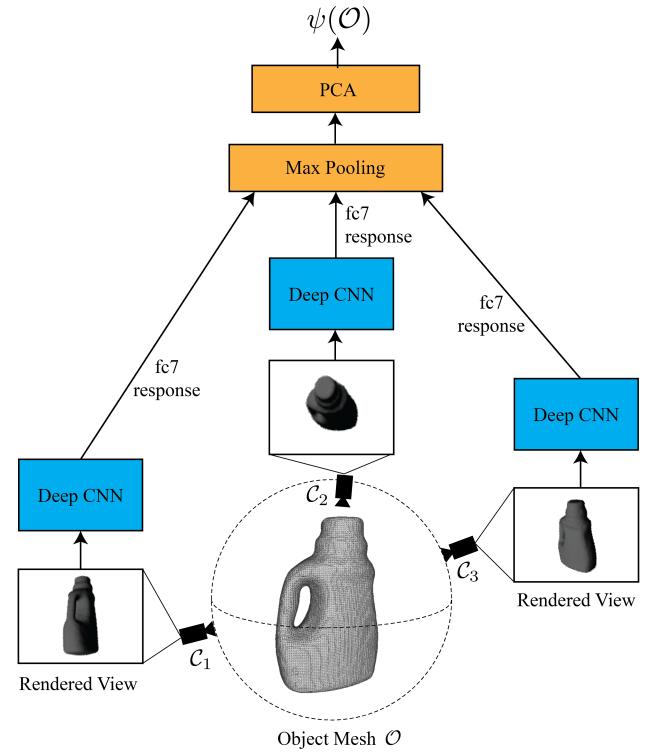


Fig. 6: Illustration of our method for embedding 3D object models in a Euclidean vector space for computing global shape similarity. We pass a set of 50 virtually rendered camera viewpoints discretized around a sphere through a deep Convolutional Neural Network (CNN) with the AlexNet [41] architecture. Finally, we take the maximum fc7 response across each of the 50 views for each dimension and run PCA to reduce the dimensionality of the output.

cross entropy loss []:

$$w_1^*, \dots, w_{N_f}^* = \underset{w_1, \dots, w_{N_f} \geq 0}{\operatorname{argmin}} \frac{1}{M} \sum_{i=1}^M \mu_i \log \left(\frac{\alpha_i}{\alpha_i + \beta_i} \right) + (1 - \mu_i) \log \left(\frac{\beta_i}{\alpha_i + \beta_i} \right)$$

where α_i and β_i are the posterior shape parameters given by Equations IV.1 and IV.2 and μ_i is the probability of force closure evaluated by exhaustive Monte-Carlo integration for a set of M grasps. We optimize this objective using SGD with the weights initialized to $w_i = 100$ for all i . [TODO: FOR EXPERIMENTS SO FAR WE ACTUALLY USE A GRID SEARCH AND NOT SGD BUT WE ARE INTERESTED IN USING THIS.]

V. GRASP SELECTION ALGORITHM

Our full algorithm for grasp selection using Multi-Armed Bandits (MABs) with prior grasp and object data from Dex-Net is detailed in Algorithm 1. We first generate a set of candidate grasps \mathcal{G} for object \mathcal{O} and predict a prior distribution for each grasp using the Dex-Net database \mathcal{D} . Next, we run MAB by selecting the next grasp using Thompson sampling, evaluating each grasp using force closure as described in Section III-D, and updating the Beta belief parameters for

each grasp. Finally, we return the grasp that has the highest mean prediction.

```

1 Input: Object  $\mathcal{O}$ , Number of Candidate Grasps  $K$ , Database
   of Prior Knowledge  $\mathcal{D}$ , Maximum Iterations  $T$ ,
2 Prior beta shape  $\alpha_0, \beta_0$ , Random Variables  $\nu, \xi$ , and  $\gamma$ 
Result: Estimate of the grasp with highest  $P_F$ ,  $\hat{\mathbf{g}}^*$ 
   // Generate candidate grasps and priors
3  $\mathcal{G} = \text{AntipodalGraspSample}(\mathcal{O}, K)$ ;
4  $\mathcal{A}_0 = \emptyset, \mathcal{B}_0 = \emptyset$ ;
5 for  $\mathbf{g}_k \in \mathcal{G}$  do
   // Equations IV.1 and IV.2
6    $\alpha_{k,0}, \beta_{k,0} = \text{ComputePriors}(\mathcal{O}, \mathbf{g}_k, \mathcal{D})$ ;
7    $\mathcal{A}_0 = \mathcal{A}_0 \cup \{\alpha_{k,0}\}, \mathcal{B}_0 = \mathcal{B}_0 \cup \{\beta_{k,0}\}$ ;
8 end
   // Run MAB with Thompson Sampling
9 for  $t = 1, \dots, T$  do
10   $j = \text{SelectNextIndex}(\mathcal{A}_{t-1}, \mathcal{B}_{t-1})$ ;
11   $\hat{\nu}, \hat{\xi}, \hat{\gamma} = \text{SampleRandomVariables}(\nu, \xi, \gamma)$ ;
12   $F_j = \text{EvaluateForceClosure}(\mathbf{g}_j, \hat{\nu}, \hat{\xi}, \hat{\gamma})$ ;
   // Equations IV.3 and IV.4
13   $\mathcal{A}_t, \mathcal{B}_t = \text{UpdateBeta}(j, F_j, \mathcal{G})$ ;
14 end
   // Find grasp with max estimated quality
15  $\mathcal{U} = \left\{ \mu_j = \frac{\alpha_{j,T}}{\alpha_{j,T} + \beta_{j,T}} \mid j = 1, \dots, K \right\}$ ;
16  $j^* = \underset{j=1, \dots, T}{\operatorname{argmax}} \mu_j$ ;
17 return  $\mathbf{g}_{j^*}$ ;
Algorithm 1: Grasp Selection With Multi-Armed Bandits
Using a Database of Prior Knowledge

```

A. Grasp Candidate Generation

In past work, Laskey et al. [44] sampled candidate grasps in 2D by sampling a grasp center from an isotropic Gaussian centered at the object center of mass and sampling a grasp angle uniformly at random. However, this method is problematic in 3D because many of the samples may be oriented away from the object and the grasps may not adequately cover the entire surface of the object.

In this work, we use a modified version of the 2D algorithm presented in Smith et al. [67] to concentrate samples on grasps that are antipodal [53]. Let K be the number of grasps to sample, w be the maximal opening of the gripper, γ be an estimate of the friction coefficient, and $\mathcal{C} = \{y \in \mathbb{R}^3 \mid |f(y)| < \epsilon\}$ be the set of points on the object surface for threshold ϵ and object SDF f . To sample a single grasp, we first generate a contact point \mathbf{c}_1 by sampling uniformly at random from \mathcal{C} , which can be done using rejection sampling. Next we form the friction cone \mathcal{F}_1 at \mathbf{c}_1 as described in Section III-C and sample a direction \mathbf{v} uniformly at random from the cone. We then compute

$$\begin{aligned} \mathbf{c}_2 &= \min_{t \geq 0} t \text{ such that } |f(\mathbf{c}_1 + (w/2 - t)\mathbf{v})| < \epsilon \\ \mathbf{x} &= 0.5(\mathbf{c}_1 + \mathbf{c}_2) \end{aligned}$$

similar to the contact computation of Section III-C and form the friction cone \mathcal{F}_2 at contact \mathbf{c}_2 . This yields a grasp $\mathbf{g} = [\mathbf{x}, \mathbf{v}]^T$. We add \mathbf{g} to our candidate set if $-\mathbf{v} \in C(\mathcal{F}_1)$ and $\mathbf{v} \in C(\mathcal{F}_2)$, where $C(\mathcal{F}) = \{y = t\mathbf{f} \mid t \in [0, \infty], \mathbf{f} \in \text{Conv}(\mathcal{F})\}$ is the convex cone for friction cone \mathcal{F} .

B. Thompson Sampling

Following Laskey et al. [44], we use Thompson sampling to select the next grasp to evaluate given belief distributions on the probability of force closure for each arm specified by estimates $\mathcal{A}_t = \{\alpha_{1,t}, \dots, \alpha_{K,t}\}$ and $\mathcal{B}_t = \{\beta_{1,t}, \dots, \beta_{K,t}\}$ at iteration t . Thompson sampling samples a probability of force closure $\hat{\theta}_{j,t} \sim B(\alpha_{j,t}, \beta_{j,t})$ for all grasps $j = 1, \dots, K$, then selects the grasp $j_t^* = \operatorname{argmax}_j \hat{\theta}_{j,t}$ to evaluate next.

Several other criteria exist for selecting the next evaluation in MAB, such as Gittins indices [44], Upper Confidence Bounds [5], and Expected Improvement [55].

VI. DEXTERITY NETWORK

[TODO: DUE TO SPACE, IT MIGHT MAKE SENSE TO DISCUSS IN EXPERIMENTS AS APPEARS TO BE COMMON PRACTICE WHEN THE DATASET IS NOT GOING TO BE RELEASED TO THE PUBLIC] It remains to describe the details of creating the Dexterity Network (Dex-Net), our prior dataset of objects and grasps with similarity, for use in the Multi-Armed Bandit models of Algorithm 1. Dex-Net 1.0 consists of approximately 20,162 3D mesh models [TODO: UPDATE WITH CORRECT COUNT] collected from the datasets described in Fig. 7. Laser-scanned models from the KIT object database [36], the Amazon Picking Challenge objects, BigBIRD [66], and YCB [7] constitute 355 of the total models. These models were chosen to reflect physical objects commonly used for benchmarking in grasping research for potential future physical experiments using Dex-Net. The majority of Dex-Net is synthetic 3D mesh models from 3DNet[73], a benchmark for object detection in robotics research, ModelNet [74], a relatively new benchmark for 3D model classification, and the SHREC 2014 large scale object retrieval challenge [49], a competition benchmark with many categories to test the state-of-the-art in shape retrieval annually. While the geometry of synthetic models does not necessarily reflect a physical object, using these models allows us to examine the effects of scale at tens of thousands of 3D models, which has not been previously attempted in grasp synthesis research.

In order to use the models for grasp selection, we first preprocess each model. Each object in Dex-Net is originally specified as a 3D mesh $\mathcal{S} = \{\mathcal{V}, \mathcal{T}\}$ where $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_V\}$ is a set of V vertices such that $\mathbf{v}_i \in \mathbb{R}^3$ for $i = 1, \dots, V$ and $\mathcal{T} = \{t_1, \dots, t_F\}$ is a set of F triangles such that $t_j \in \mathbb{Z}_+^3$ for $j = 1, \dots, F$. First, we remove unreferenced vertices and illegal triangles (e.g. same vertex is referenced twice in one triangle). Next, we normalize the orientation of the model by performing Principal Component Analysis (PCA) on the raw vertices of the mesh and then rotating the object such that its first principal component aligns with the z -axis and its second aligns with the y -axis. We then set the object center of mass \mathbf{z} to be the center of the bounding box for the reoriented mesh. Since synthetic models may not be specified in meters, we also rescale each mesh such that the smallest dimension of the bounding box lies within the maximum opening of the gripper. Finally, we convert each mesh to

Dataset	# Models	Synthetic?	# Categories
KIT Object Database	129	N	N/A
Amazon Picking Challenge	26	N	N/A
BigBIRD	120	N	N/A
YCB	80	N	N/A
3D Net	1371	Y	50
SHREC 2014 Large Scale Challenge	8987	Y	171
ModelNet	9449	Y	40

Fig. 7: The seven datasets used in Dex-Net 1.0 with information on whether or not the models were laser-scanned, the number of models in the dataset, and the number of labelled categories (if category labels are present). [TODO: UPDATE MODELNET NUMBER - WE ACTUALLY ONLY USE SELECT CATEGORIES]

a signed distance field f using SDFGen [], an open-source C++ tool.

Once each 3D model is preprocessed, we compute the shape feature vector described in Section IV-B.3 for each object. Each feature vector is then put into a KD-Tree nearest neighbor query structure for accelerated lookups of shape nearest neighbors. This structure consistsutes our shape similiarity network, in which each object \mathcal{O} is connected to its n nearest neighbors $\mathcal{O}_1, \dots, \mathcal{O}_n$ by an edge of weight $\|\psi(\mathcal{O}) - \psi(\mathcal{O}_i)\|_2^2$. We generate a set of grasps for each object in the network using the method of Section V-A and evaluate the probability of force closure for each grasp using brute-force Monte Carlo integration as a benchmark [38]. As the number of models is quite large, we distribute the grasp labelling for each object across virtual machines in Google Compute Engine and aggregate the results at the end.

VII. EXPERIMENTS

[TODO: ALL EXPERIMENTS HERE ARE PRELIMINARY. ADD NEW RESULTS FROM NON-SYMMETRIC OBJECTS AND MORE DATA] [TODO: INCLUDE EXPERIMENTS ON SENSITIVITY TO UNCERTAINTY]
[TODO: IMAGES OF SELECTED GRASPS]

We evaluated the convergence rate of Algorithm 1 for varying sizes of prior data used from Dex-Net and examined the sensitivity to parameters of the algorithm. To test the generalization of our feature maps, we reserved objects from YCB and ModelNet as a test set, and used nested subsets of the remaining data as training data for optimizing the hyperparameters and for tuning the MAB convergence rate. We assumed a PR2 gripper for our grasps, which has a maximum width of approximately 10cm. We assumed isotropic covariance with object and gripper translation variance $\sigma_t = 0.005$, object and gripper rotation variance $\sigma_r = 0.1$, and friction variance $\sigma_\gamma = 0.1$. The weights of our similarity kernel were selected using a grid search over possible values on a held-out validation set.

A. System

To handle the scale of the experiments, we developed a Cloud-based software library on top of Google Cloud

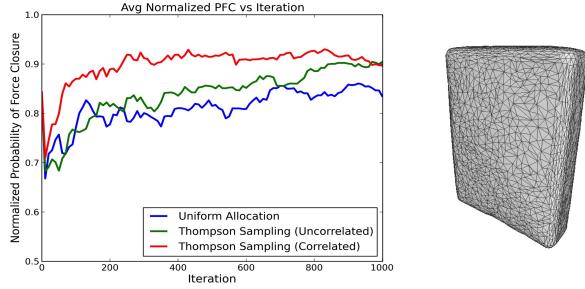
services. We used Google Cloud Storage (GCS) to store the mesh models and to sample and evaluate 250 grasps per object. We used Google Compute Engine to distribute trials of MAB algorithms across objects, reducing the runtime of our large-scale experiments. Our master analysis script can launch up to 600 GCE single core instances, each running Ubuntu 12.04. When an instance is launched, we mounted Dex-Net via a persistent disk containing a particular set of objects and grasps. Each instance pulls our latest code from Github, runs a single python experiment script configured in a YAML file, compresses the results in an output directory, and uploads the results to GCS. After all instances complete, the master script turns off all instances, unmounts all disks, downloads the results from each instance, and sends the user a notification email. The results can be optionally analyzed by the master script.

B. Rate of Convergence

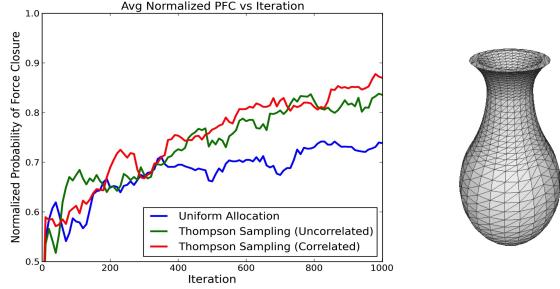
[TODO: REPLACE WITH FINAL VERSIONS. THE CURRENT RESULTS ARE PRELIMINARY]

Fig. 8 shows the normalized P_F (the ratio of the P_F for the sampled grasp to the P_F in the candidate grasp set) versus iteration averaged over 20 trials for a cereal box and flower pot. The plot compares Algorithm 1 with no prior data to uncorrelated Thompson sampling and uniform allocation, two methods used by Laskey et al. [44]. We see that the CCBP model outperforms the uncorrelated model by approximately 10× for the cereal box but has comparable performance to Thompson sampling on the flower pot. Iniital reuslts suggest that performance on the flower pot is poorer because few grasps are similar according to our feature representation.

To examine the effects of orders of magnitude of prior data in the MAB algorithms, we ran the algorithms with priors computed from increasingly larger subsets of prior data from Dex-Net: 0, 15, 150, and 1500 objects. Fig. 9 shows the normalized probability of force closure P_F versus iteration when planning grasps for a bottle and detergent averaged over 50 trials. For each run of our algorithm, we used prior grasp from the five closest objects in Dex-Net in the multi-view CNN-based shape similiarity vector space. The objects from Dex-Net were labelled with 250 grasps and P_F evaluated using brute force Monte-Carlo integration. For the bottle, we see that the convergence of the correlated MAB algorithms to a grasp with high P_F accelerates with increasingly larger subsets of prior data used, and is approximately 10× faster when using 1500 prior objects. For the detergent the gains are more modest, with the largest dataset accelerating convergence to within 90% of the optimal grasp in the set by approximately 1.5×. As illustrated in Fig. 10, the nearest neighbor objects become increasingly similar to the object to label with increasing sizes of datasets from Dex-Net. This suggests that more prior data may also lead to more similar objects to rarer categories such as the detergent bottle, which may accelerate convergence [TODO: REMOVE] Future work will examine



(a) Normalized P_F on a cereal box with 250 candidate parallel-jaw grasps.



(b) Normalized P_F on a flower pot with 250 candidate parallel-jaw grasps.

Fig. 8: Comparison of the normalized P_F of the sampled parallel-jaw grasp versus iteration of the MAB algorithm for correlated Thompson sampling, uncorrelated Thompson sampling, and Uniform Allocation averaged over 20 trials. (a) On the cereal box, correlated Thompson sampling converges to within 90% of the highest quality grasp approximately 4× faster than uncorrelated. (b) However, correlated and uncorrelated Thompson sampling perform comparably on the flower pot because there are few grasp similarities in the candidate set.

these effects on a larger subset of objects and will use additional nearest neighbors from Dex-Net to compute priors.

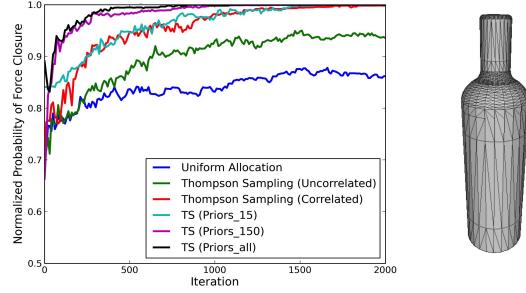
VIII. SENSITIVITY TO UNCERTAINTY

[TODO: WRITE SECTION]

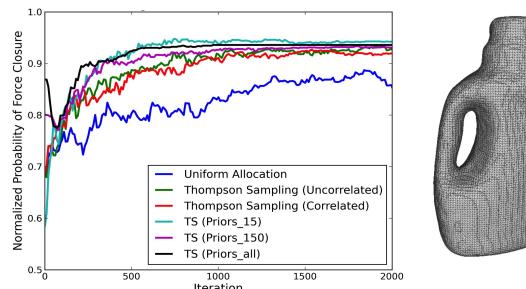
IX. DISCUSSION AND FUTURE WORK

We introduced the Dexterity Network (Dex-Net) 1.0, a dataset of approximately 20,000 3D object models and a network of similarity between object models and parallel-jaw grasps on each model. We extended the Multi-Armed Bandit model of Laskey et al. [44] to leverage similarity between grasps based on gripper pose, local surface geometry encoded as heightmaps, and global shape similarity encoded using a state-of-the-art method for shape retrieval based on multi-view Convolutional Neural Networks (CNNs) [70]. Our experiments suggest that increasing the amount of prior knowledge of objects and the quality of grasps on each object enables up to a 10× speedup in the time to select a grasp with high probability of force closure for a new object, and also suggest that convergence is fastest when the object has a set of geometrically similar nearest neighbors in the database.

One current shortcoming of our method is that it relies on a set of features manually designed to predict grasp quality for a grasp and object, such as heightmaps and CNNs that predict object category. Future work will leverage recent developments in feature optimization with deep neural networks [48] to optimize feature representations of both the local and global surface geometry based on the outcome of



(a) Normalized P_F on a bottle with 250 candidate parallel-jaw grasps.



(b) Normalized P_F on a detergent bottle with 250 candidate parallel-jaw grasps.

Fig. 9: Comparison of the normalized P_F of the sampled parallel-jaw grasp (y-axis) versus iteration of the MAB algorithm (x-axis) for correlated Thompson sampling with and without the use of prior grasps from five nearest neighbors from Dex-Net, uncorrelated Thompson sampling, and Uniform Allocation averaged over 20. Prior grasps were taken from increasingly larger subsets of Dex-Net: 15, 150, and 1500 objects. For the bottle, convergence to within 90% of the optimal grasp is accelerated by approximately 10× over uncorrelated Thompson sampling, and performance appears to improve with increasing sizes of the prior dataset used. The detergent converges slightly faster with more prior data, however the speedup is approximately 1.5× and does not appear to improve with additional data from Dex-Net.

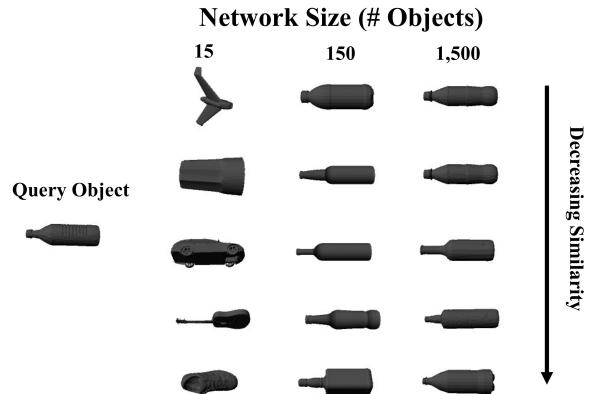


Fig. 10: Illustration of the five nearest neighbor objects from Dex-Net using the multi-view CNN-based shape embedding and Euclidean distance for increasing sizes of randomly selected subsets of data from Dex-Net. We see that while the neighbors for a small dataset are not relevant, the neighbors become increasingly similar to the query bottle as more data is used.

grasp evaluations in the bandit algorithms. We will also study extensions of the shape similarity neural network to predict object class from simulated depth images of objects on a table, which could be useful for transferring learned data from Dex-Net to objects in the physical world.

Another shortcoming is our current use of wrench space metrics, which have been criticized recently for not correlating with human labels [35] or physical trials [1]. However,

both forms of labels may be difficult to acquire for the entire dataset because of the human and time cost involved. Future work will research using Multi-Armed Bandits (MABs) to actively acquire grasp evaluations uncertain objects. To do so, we will study using functional map networks to establish consistent maps between prototypical object parts that indicate grasp affordance, such as handles. We will also build on recent developments in caging to label key the dataset with caging hand configurations either analytically, using simulation, or using a learned model to predict caging configurations. We will also work toward releasing Dex-Net as an open-access project with an open-source API for labelling objects in the Cloud with analytic quality metrics or simulation outcomes, integrating Dex-Net with physical robots, and learning statistical models of grasp success using deep learning and MAB.

REFERENCES

- [1] R. Balasubramanian, L. Xu, P. D. Brook, J. R. Smith, and Y. Matsuoka, “Physical human interactive guidance: Identifying grasping principles from human-planned grasps,” *Robotics, IEEE Transactions on*, vol. 28, no. 4, pp. 899–910, 2012.
- [2] T. D. Barfoot and P. T. Furgale, “Associating uncertainty with three-dimensional poses for use in estimation problems,” *Robotics, IEEE Transactions on*, vol. 30, no. 3, pp. 679–693, 2014.
- [3] J. Bohg, A. Morales, T. Asfour, and D. Kragic, “Data-driven grasp synthesis survey,” *Robotics, IEEE Transactions on*, vol. 30, no. 2, pp. 289–309, 2014.
- [4] A. Boularias, J. A. Bagnell, and A. Stentz, “Efficient optimization for autonomous robotic manipulation of natural objects,” 2014.
- [5] ———, “Learning to manipulate unknown objects in clutter by reinforcement,” in *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [6] P. Brook, M. Ciocarlie, and K. Hsiao, “Collaborative grasp planning with multiple object representations,” in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*. IEEE, 2011, pp. 2851–2858.
- [7] B. Calli, A. Walsman, A. Singh, S. Srinivasa, P. Abbeel, and A. M. Dollar, “Benchmarking in manipulation research: The ycb object and model set and benchmarking protocols,” *arXiv preprint arXiv:1502.03143*, 2015.
- [8] J.-S. Cheong, H. Kruger, and A. F. van der Stappen, “Output-sensitive computation of force-closure grasps of a semi-algebraic object,” *IEEE Trans. Automation Science and Engineering*, vol. 8, no. 3, pp. 495–505, 2011.
- [9] 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.
- [10] W. Chu, L. Li, L. Reyzin, and R. E. Schapire, “Contextual bandits with linear payoff functions,” in *International Conference on Artificial Intelligence and Statistics*, 2011, pp. 208–214.
- [11] C. Cieri, D. Miller, and K. Walker, “The fisher corpus: a resource for the next generations of speech-to-text,” in *LREC*, vol. 4, 2004, pp. 69–71.
- [12] M. T. Ciocarlie and P. K. Allen, “Hand posture subspaces for dexterous robotic grasping,” *Int. J. Robotics Research (IJRR)*, vol. 28, no. 7, pp. 851–867, 2009.
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 2009, pp. 248–255.
- [14] R. Detry, C. H. Ek, M. Madry, and D. Kragic, “Learning a dictionary of prototypical grasp-predicting parts from grasping experience,” in *Robotics and Automation (ICRA), 2013 IEEE International Conference on*. IEEE, 2013, pp. 601–608.
- [15] R. Detry, C. H. Ek, M. Madry, J. Piater, and D. Kragic, “Generalizing grasps across partly similar objects,” in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 3791–3797.
- [16] R. Detry, D. Kraft, O. Kroemer, L. Bodenhausen, J. Peters, N. Krüger, and J. Piater, “Learning grasp affordance densities,” *Paladyn, Journal of Behavioral Robotics*, vol. 2, no. 1, pp. 1–17, 2011.
- [17] R. Diankov, “Automated construction of robotic manipulation programs,” Ph.D. dissertation, Citeseer, 2010.
- [18] S. Dragiev, M. Toussaint, and M. Gienger, “Uncertainty aware grasping and tactile exploration,” in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*. IEEE, 2013, pp. 113–119.
- [19] C. Ferrari and J. Canny, “Planning optimal grasps,” in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 1992, pp. 2290–2295.
- [20] D. Fischinger, A. Weiss, and M. Vincze, “Learning grasps with topographic features,” *Int. J. Robotics Research (IJRR)*, p. 0278364915577105, 2015.
- [21] R. Goetschalckx, P. Poupart, and J. Hoey, “Continuous correlated beta processes,” in *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*, vol. 22, no. 1. Citeseer, 2011, p. 1269.
- [22] 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.
- [23] C. Goldfeder and P. K. Allen, “Data-driven grasping,” *Autonomous Robots*, vol. 31, no. 1, pp. 1–20, 2011.
- [24] C. Goldfeder, M. Ciocarlie, H. Dang, and P. K. Allen, “The columbia grasp database,” in *Robotics and Automation, 2009. ICRA'09. IEEE International Conference on*. IEEE, 2009, pp. 1710–1716.
- [25] C. Goldfeder, M. Ciocarlie, J. Peretzman, H. Dang, and P. K. Allen, “Data-driven grasping with partial sensor data,” in *Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on*. IEEE, 2009, pp. 1278–1283.
- [26] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Satheesh, S. Sengupta, A. Coates, et al., “Deep-speech: Scaling up end-to-end speech recognition,” *arXiv preprint arXiv:1412.5567*, 2014.
- [27] J. Hays and A. A. Efros, “Im2gps: estimating geographic information from a single image,” in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*. IEEE, 2008, pp. 1–8.
- [28] A. Herzog, P. Pastor, M. Kalakrishnan, L. Righetti, T. Asfour, and S. Schaal, “Template-based learning of grasp selection,” in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 2379–2384.
- [29] A. Herzog, P. Pastor, M. Kalakrishnan, L. Righetti, J. Bohg, T. Asfour, and S. Schaal, “Learning of grasp selection based on shape-templates,” *Autonomous Robots*, vol. 36, no. 1-2, pp. 51–65, 2014.
- [30] U. Hillenbrand, M. Roa, et al., “Transferring functional grasps through contact warping and local replanning,” in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*. IEEE, 2012, pp. 2963–2970.
- [31] K. Hsiao, M. Ciocarlie, and P. Brook, “Bayesian grasp planning,” in *ICRA 2011 Workshop on Mobile Manipulation: Integrating Perception and Manipulation*, 2011.
- [32] K. Hsiao, L. P. Kaelbling, and T. Lozano-Perez, “Grasping pomdps,” in *Robotics and Automation, 2007 IEEE International Conference on*. IEEE, 2007, pp. 4685–4692.
- [33] Y. Jiang, S. Mosenes, and A. Saxena, “Efficient grasping from rgbd images: Learning using a new rectangle representation,” in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE, 2011, pp. 3304–3311.
- [34] G. Kahn, P. Sujan, S. Patil, S. Bopardikar, J. Ryde, K. Goldberg, and P. Abbeel, “Active exploration using trajectory optimization for robotic grasping in the presence of occlusions.”
- [35] D. Kappler, J. Bohg, and S. Schaal, “Leveraging big data for grasp planning,” in *Robotics and Automation (ICRA), 2015 IEEE International Conference on*. IEEE, 2015, pp. 4304–4311.
- [36] A. Kasper, Z. Xue, and R. Dillmann, “The kit object models database: An object model database for object recognition, localization and manipulation in service robotics,” *The International Journal of Robotics Research*, vol. 31, no. 8, pp. 927–934, 2012.
- [37] B. Kehoe, D. Berenson, and K. Goldberg, “Estimating part tolerance bounds based on adaptive cloud-based grasp planning with slip,” in *Proc. IEEE Conf. on Automation Science and Engineering (CASE)*. IEEE, 2012, pp. 1106–1113.
- [38] ———, “Toward cloud-based grasping with uncertainty in shape: Estimating lower bounds on achieving force closure with zero-slip push grasps,” in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*. IEEE, 2012, pp. 576–583.

- [39] B. Kehoe, A. Matsukawa, S. Candido, J. Kuffner, and K. Goldberg, “Cloud-based robot grasping with the google object recognition engine,” in *Robotics and Automation (ICRA), 2013 IEEE International Conference on*. IEEE, 2013, pp. 4263–4270.
- [40] J. Kim, K. Iwamoto, J. J. Kuffner, Y. Ota, and N. S. Pollard, “Physically-based grasp quality evaluation under uncertainty,” in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*. IEEE, 2012, pp. 3258–3263.
- [41] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [42] O. Kroemer, R. Detry, J. Piater, and J. Peters, “Combining active learning and reactive control for robot grasping,” *Robotics and Autonomous Systems*, vol. 58, no. 9, pp. 1105–1116, 2010.
- [43] H. Kruger *et al.*, “Partial closure grasps: Metrics and computation,” in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE, 2011, pp. 5024–5030.
- [44] M. Laskey, J. Mahler, Z. McCarthy, F. Pokorny, S. Patil, J. van den Berg, D. Kragic, P. Abbeel, and K. Goldberg, “Multi-arm bandit models for 2d sample based grasp planning with uncertainty.” in *Proc. IEEE Conf. on Automation Science and Engineering (CASE)*. IEEE, 2015.
- [45] I. Lenz, R. Knepper, and A. Saxena, “Deepmpc: Learning deep latent features for model predictive control,” in *Proc. Robotics: Science and Systems (RSS)*, 2015.
- [46] I. Lenz, H. Lee, and A. Saxena, “Deep learning for detecting robotic grasps,” *The International Journal of Robotics Research*, vol. 34, no. 4-5, pp. 705–724, 2015.
- [47] B. León, S. Ulbrich, R. Diankov, G. Puche, M. Przybylski, A. Morales, T. Asfour, S. Moisio, J. Bohg, J. Kuffner, *et al.*, “Opengrasp: a toolkit for robot grasping simulation,” in *Simulation, Modeling, and Programming for Autonomous Robots*. Springer, 2010, pp. 109–120.
- [48] S. Levine, C. Finn, T. Darrell, and P. Abbeel, “End-to-end training of deep visuomotor policies,” *arXiv preprint arXiv:1504.00702*, 2015.
- [49] B. Li, Y. Lu, C. Li, A. Godil, T. Schreck, M. Aono, M. Burtscher, Q. Chen, N. K. Chowdhury, B. Fang, *et al.*, “A comparison of 3d shape retrieval methods based on a large-scale benchmark supporting multimodal queries,” *Computer Vision and Image Understanding*, vol. 131, pp. 1–27, 2015.
- [50] Y. Li and N. S. Pollard, “A shape matching algorithm for synthesizing humanlike enveloping grasps,” in *Humanoid Robots, 2005 5th IEEE-RAS International Conference on*. IEEE, 2005, pp. 442–449.
- [51] 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.
- [52] S. Liu and S. Carpin, “A fast algorithm for grasp quality evaluation using the object wrench space,” 2015.
- [53] 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,” 2015.
- [54] A. T. Miller and P. K. Allen, “Graspit! a versatile simulator for robotic grasping,” *Robotics & Automation Magazine, IEEE*, vol. 11, no. 4, pp. 110–122, 2004.
- [55] L. Montesano and M. Lopes, “Active learning of visual descriptors for grasping using non-parametric smoothed beta distributions,” *Robotics and Autonomous Systems*, vol. 60, no. 3, pp. 452–462, 2012.
- [56] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, “Kinectfusion: Real-time dense surface mapping and tracking,” in *Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on*. IEEE, 2011, pp. 127–136.
- [57] F. Panahi, M. Davoodi, and A. F. van der Stappen, “Orienting parts with shape variation,” in *Algorithmic Foundations of Robotics XI*. Springer, 2015, pp. 479–496.
- [58] F. T. Pokorny, Y. Bekiroglu, and D. Kragic, “Grasp moduli spaces and spherical harmonics,” in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE, 2014, pp. 389–396.
- [59] F. T. Pokorny, K. Hang, and D. Kragic, “Grasp moduli spaces.” in *Robotics: Science and Systems*, 2013.
- [60] F. T. Pokorny and D. Kragic, “Classical grasp quality evaluation: New algorithms and theory,” in *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*. IEEE, 2013, pp. 3493–3500.
- [61] N. S. Pollard, “Parallel methods for synthesizing whole-hand grasps from generalized prototypes,” DTIC Document, Tech. Rep., 1994.
- [62] D. Prattichizzo and J. C. Trinkle, “Grasping,” in *Springer handbook of robotics*. Springer, 2008, pp. 671–700.
- [63] J. Redmon and A. Angelova, “Real-time grasp detection using convolutional neural networks,” *arXiv preprint arXiv:1412.3128*, 2014.
- [64] M. Salganicoff, L. H. Ungar, and R. Bajcsy, “Active learning for vision-based robot grasping,” *Machine Learning*, vol. 23, no. 2-3, pp. 251–278, 1996.
- [65] A. Saxena, J. Driemeyer, and A. Y. Ng, “Robotic grasping of novel objects using vision,” *The International Journal of Robotics Research*, vol. 27, no. 2, pp. 157–173, 2008.
- [66] A. Singh, J. Sha, K. S. Narayan, T. Achim, and P. Abbeel, “Bigbird: A large-scale 3d database of object instances,” in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE, 2014, pp. 509–516.
- [67] G. Smith, E. Lee, K. Goldberg, K. Bohringer, and J. Craig, “Computing parallel-jaw grips,” in *Robotics and Automation, 1999. Proceedings. 1999 IEEE International Conference on*, vol. 3. IEEE, 1999, pp. 1897–1903.
- [68] T. Stouraitis, U. Hillenbrand, and M. A. Roa, “Functional power grasps transferred through warping and replanning,” in *Robotics and Automation (ICRA), 2015 IEEE International Conference on*. IEEE, 2015, pp. 4933–4940.
- [69] 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.
- [70] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, “Multi-view convolutional neural networks for 3d shape recognition,” *arXiv preprint arXiv:1505.00880*, 2015.
- [71] F. Tombari, S. Salti, and L. Di Stefano, “Unique signatures of histograms for local surface.”
- [72] 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.
- [73] W. Wohlkinger, A. Aldoma, R. B. Rusu, and M. Vincze, “3dnet: Large-scale object class recognition from cad models,” in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 5384–5391.
- [74] Z. Wu, S. Song, A. Khosla, X. Tang, and J. Xiao, “3d shapenets for 2.5 d object recognition and next-best-view prediction,” *arXiv preprint arXiv:1406.5670*, 2014.
- [75] 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.