# Dex-Net 1.0: A Cloud-Based Network of 3D Objects for Robust Grasp Planning Using a Multi-Armed Bandit Model with Correlated Rewards

Jeffrey Mahler[1], Florian T. Pokorny[1], Brian Hou[1], Melrose Roderick[1], Michael Laskey[1], Mathieu Aubry[1,2], Kai Kohlhoff[3], Torsten Kroeger[3], James Kuffner[3], Ken Goldberg[1]

*Abstract*— This paper presents Dexterity Network 1.0 (Dex-Net), a new dataset and associated algorithm to study the scaling effects of Big Data and Cloud Computation on robust grasp planning. The algorithm uses a Multi-Armed Bandit model with correlated rewards to leverage prior grasps and 3D object models in a growing dataset that currently includes over 10,000 unique 3D object models and 2.5 million parallel-jaw grasps. Each grasp includes an estimate of the probability of force closure under uncertainty in object and gripper pose and friction. Dex-Net 1.0 uses Multi-View Convolutional Neural Networks (MV-CNNs), a new deep learning method for 3D object classification, as a similarity metric between objects and the Google Cloud Platform to simultaneously run up to 1,500 virtual cores, reducing runtime by three orders of magnitude. Experiments suggest that using prior data can significantly benefit the quality and complexity of robust grasp planning. We report on system sensitivity to varying similarity metrics and pose and friction uncertainty levels. Code and additional information can be found at: **http://berkeleyautomation. github.io/dex-net/**.

## I. INTRODUCTION

Cloud-based Robotics and Automation systems exchange data and perform computation via networks instead of operating in isolation with limited computation and memory. Potential advantages to using the Cloud include Big Data: access to updated libraries of images, maps, and object/product data; and Parallel Computation: access to parallel grid computing for statistical analysis, machine learning, and planning [22]. These benefits have recently been demonstrated in vision and speech, where datasets with millions of examples such as ImageNet have produced results [14], [24] that surpass those obtained from decades of research on analytic methods. This suggests that large-scale machine learning of grasps for vast numbers of possible object shapes, object poses, and environment configurations [12], [18], [27], could exhibit scaling effects similar to those observed in computer vision and speech recognition.

The primary contribution of this paper is an algorithm based on a Multi-Armed Bandit (MAB) model with correlated rewards to speed up robust planning by learning from a large dataset of prior grasps and 3D object models. The algorithm is based on Continuous Correlated Beta Processes (CCBPs) [11], [32], an efficient model for predicting a belief distribution on the quality of each grasp from prior data.

[1] University of California, Berkeley, USA; {jmahler, ftpokorny, brian.hou, goldberg}@berkeley.edu, melrose_roderick@brown.edu

[2] Université Paris-Est, LIGM (UMR CNRS 8049), ENPC, France. mathieu.aubry@enpc.fc

[3] Google Inc., Mountain View, USA; {kohlhoff, tkr, kuffner}@google.com
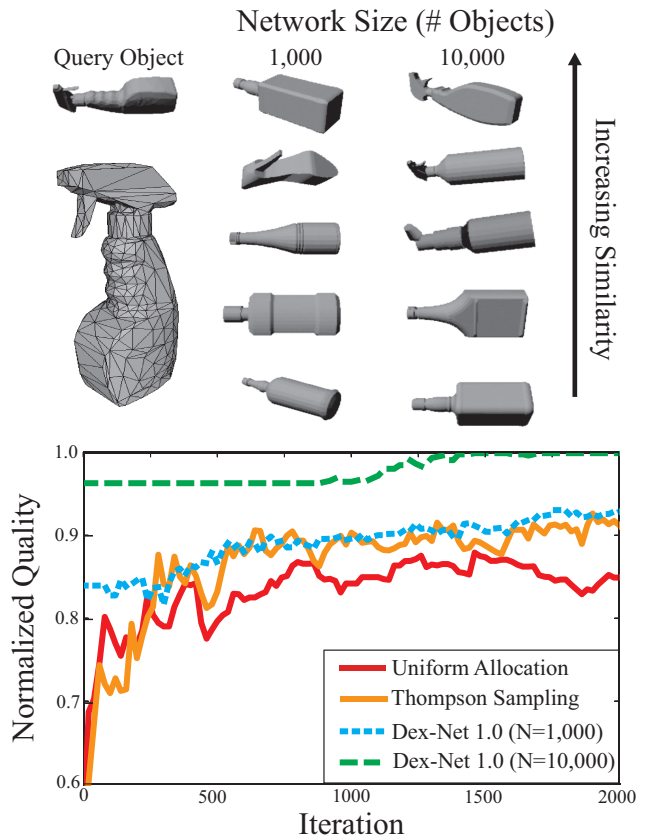
Fig. 1: Normalized grasp quality versus iteration averaged over 25 trials for the Dex-Net 1.0 algorithm with 1,000 and 10,000 prior 3D objects (bottom) and illustrations of five nearest neighbors in Dex-Net (top) for a spray bottle. We measure quality by the probability of force closure of the best grasp predicted by the algorithm on each iteration and compare with Thompson sampling without priors and uniform allocation. (Top) The spray bottle has no similar neighbors with 1,000 objects, but two other spray bottles are found by the MV-CNN in the 10,000 object set. (Bottom) As a result, the Dex-Net 1.0 algorithm does not outperform Thompson sampling for 1,000 objects, but quickly converges to the optimal grasp with 10,000 prior objects.

To study scaling effects, we developed Dex-Net 1.0, a growing dataset that currently includes over 10,000 unique 3D object models selected to reflect objects that could be encountered in warehousing or the home such as containers, tools, tableware, and toys. Dex-Net also contains approximately 2.5 million parallel-jaw grasps, as each object is labelled with up to 250 grasps and an estimate of the probability of force closure for each under uncertainty in object pose, gripper pose, and friction coefficient. To the best of our knowledge, this is the largest object dataset used for grasping research to-date. We also incorporate Multi-View Convolutional Neural Networks (MV-CNNs) [43], a state-of-the-art method for 3D shape classification, to efficiently retrieve similar 3D objects.

We implement the algorithm on Google Compute Engine and store Dex-Net 1.0 on Google Cloud Storage, with a system that can run up to 1,500 instances at once. Experiments suggest that using 10,000 prior object models from Dex-Net reduces the number of samples needed to plan robust parallel-jaw grasps by up to $2\times$ on average over 45 objects.

## II. RELATED WORK

Grasp planning considers the problem of finding grasps for a given object that achieve force closure or optimize a related quality metric [36]. Usually it is assumed that the object is known exactly and that contacts are placed exactly, and mechanical wrench space analysis is applied. Robust grasp planning considers the same problem in the presence of bounded perturbations in properties such as object shape, pose, or mechanical properties such as friction, which are inevitable due to imprecision in perception and control. One way to treat perturbations is statistical sampling. Since sampling in high dimensions can be computationally demanding, recent work has explored how a robust grasp computed from one object can guide the search for robust grasps on similar objects [4], for example by warping contacts [42] or interpolating grasps and shapes over a Grasp Moduli Space [35]. To study grasp planning at scale, Goldfeder et al. [12], [13] developed the Columbia grasp database, a dataset of 1,814 distinct models and over 200,000 force closure grasps generated using the GraspIt! stochastic sampling-based grasp planner.

Recent research has studied labelling grasps in a database with metrics that are robust to imprecision in perception and control using probability of force closure ($P_F$) [44] or expected Ferrari-Canny quality [23]. Experiments by Weisz et al. [44] and Kim et al. [23] suggest that the robust metrics are better correlated with success on a physical robot than deterministic wrench space metrics. Brook et al. [6] planned robust grasps for a database of 892 point clouds and developed a model to predict grasp success on a physical robot based on correlations with grasps in the database. Kehoe et al. [21] created a Cloud-based system to transfer grasps evaluated by $P_F$ on 100 objects in a database to a physical robot by indexing the objects with the Google Goggles object recognition engine, and achieved 80% success in grasping objects on a table.

Another line of research has focused on synthesizing grasps using statistical models [4] learned from a database of images [27] or point clouds [9], [15], [47] of objects annotated with grasps from human demonstrators [15], [27] or physical execution [15]. Kappler et al. [18] created a database of over 700 object instances, each labelled with 500 Barrett hand grasps and their associated quality from human annotations and the results of simulations with the ODE physics engine. The authors trained a deep neural network to predict grasp quality from heightmaps of the local object surface. We estimate $P_F$ using similar objects and grasps using a variant of the Multi-Armed Bandit (MAB) model for sequential decision-making.

Our work is also closely related to research on actively sampling grasps to build a statistical model of grasp quality from fewer examples [10], [25], [37]. Montesano and Lopes [32] used Continuous Correlated Beta Processes [11] to actively acquire grasp executions on a physical robot, and measured correlations from the responses to a bank of image filters designed to detect grasp affordances such as edges. Oberlin and Tellex [33] developed a budgeted MAB algorithm for planning 3 DOF crane grasps using priors from the responses of hand-designed depth image filters, but did not study the effects of orders of magnitude of prior data on convergence. Recently, Laskey et al. [26] used MAB algorithms to reduce the number of samples needed to identify grasps with high $P_F$ under uncertainty in object shape, pose, and friction in 2D. In this work we extend the model of [26] to 3D and study the scaling effects of using prior data from Dex-Net on planning grasps with high $P_F$.

To use the prior information contained in Dex-Net, we also draw on research on 3D model similarity. One line of research has focused on shape geometry, such as characteristics of harmonic functions on the shape [5], or CNNs trained on a voxel representation of shape [31], [46]. Another line of research relies on the description of rendered views of a 3D model [8], [12]. One of the key difficulty of these methods is comparing views from different objects, which may be oriented inconsistently. The recent work of Su et al. [43] addresses this issue by using CNN trained for ImageNet classification as descriptors for the different views and aggregating them with a second CNN that learns the invariance to orientation. Using this method, the authors improve state-of-the-art classification accuracy on ModelNet40 by 10%. We use a max-pooling to aggregate views, similar to the average pooling proposed in [1].

## III. DEFINITIONS AND PROBLEM STATEMENT

We consider the robust grasp planning problem for a given 3D object model and parallel-jaw grippers using probability of force closure ($P_F$) under uncertainty in object pose, gripper pose, and friction coefficient as a grasp quality metric. We assume the nominal object shape is given as a signed distance function (SDF) $f : \mathbb{R}^3 \to \mathbb{R}$ [30], which is zero on the object surface, positive outside the object, and negative within. The object is specified in units of meters with given center of mass $\mathbf{z} \in \mathbb{R}^3$. We assume soft-finger contacts with a Coulomb friction model [48] and that the gripper jaws are opened to their maximal width $w \in \mathbb{R}$ before closing on the object.

### A. Grasp and Object Parameterization

The grasp parameters are illustrated in Fig. 2. Let $\mathbf{g} = (\mathbf{x}, \mathbf{v})$ 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 denote by $\mathcal{S} = \{\mathbf{y} \in \mathbb{R}^3 \big| f(\mathbf{y}) = 0\}$ the surface of an object for SDF $f$, and specify all points with respect to a reference frame centered at the object center of mass $\mathbf{z}$ and oriented along the principal axes of $\mathcal{S}$. Let $\mathcal{G} = \{(\mathbf{x}, \mathbf{v}) \big| \mathbf{x} \in \mathbb{R}^3, \mathbf{v} \in \mathbb{S}^2\}$ denote the space of all grasps and $\mathcal{H} = \{\mathcal{O} = \{\mathbf{z}, f(\cdot)\} \big| \mathbf{z} \in \mathbb{R}^3, f \in \mathcal{A}\}$ denote the space of all objects, where $\mathcal{A}$ is the space of all SDFs for closed
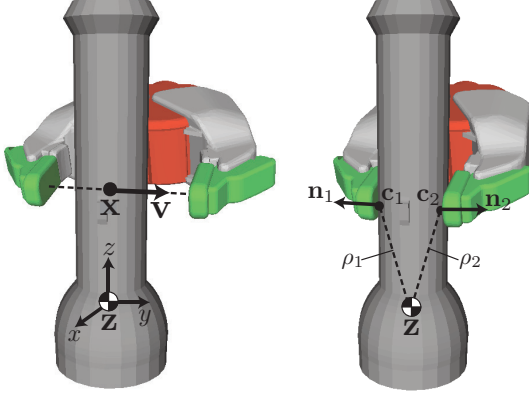
Fig. 2: Illustration of grasp parameterization and contact model. (Left) We parameterize parallel-jaw grasps by the centroid of the jaws $\mathbf{x} \in \mathbb{R}^3$ and approach direction, or direction along which the jaws close, $\mathbf{v} \in \mathbb{S}^2$. The parameters $\mathbf{x}$ and $\mathbf{v}$ are specified with respect to a coordinate frame 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 $\rho_i = \mathbf{c}_i - \mathbf{z}$.

and compact surfaces. We denote by $\mathcal{M} = \mathcal{G} \times \mathcal{H}$ the Grasp Moduli Space of all parallel-jaw grasps and objects [35].

### B. Sources of Uncertainty

We assume Gaussian distributions on object pose, gripper pose, and friction coefficient to model errors in registration, robot calibration, or classification material properties, respectively. We denote by $\mathcal{N}(\mu, \Sigma)$ a Gaussian distribution with mean $\mu$ and variance $\Sigma$. Let $\upsilon \sim \mathcal{N}(\mathbf{0}, \Sigma_\upsilon)$ denote a zero-mean Gaussian on $\mathbb{R}^6$ and $\mu_\xi \in SE(3)$ be a mean object pose. We define the object pose random variable $\xi = \exp(\upsilon^\wedge)\mu_\xi$, where the $\wedge$ operator maps from $\mathbb{R}^6$ to the Lie algebra $\mathfrak{se}(3)$ [2]. Let $\nu \sim \mathcal{N}(\mathbf{0}, \Sigma_\nu)$ denote zero-mean Gaussian gripper pose uncertainty with mean $\mu_\nu \in \mathcal{G}$. Let $\gamma \sim \mathcal{N}(\mu_\gamma, \Sigma_\gamma)$ denote a Gaussian distribution on friction coefficient with mean $\mu_\gamma \in \mathbb{R}$. We denote by $\hat{\xi}$, $\hat{\nu}$, and $\hat{\gamma}$ samples of the random variables.

### C. Contact Model

Given a grasp $\mathbf{g}$ on an object $\mathcal{O}$ and samples $\hat{\xi}, \hat{\nu}$, and $\hat{\gamma}$, let $\mathbf{c}_i \in \mathbb{R}^3$ for $i \in 1, 2$ denote the 3D contact location between the $i$-th gripper jaw and surface as shown in Fig. 2. Each contact $\mathbf{c}_i = \mathbf{x} + (-1)^i(w/2 - t_i^*)\mathbf{v}$ where [30]:

$$t_i^* = \underset{t \geqslant 0}{\operatorname{argmin}}\ t \text{ such that } f\left(\mathbf{x} + (-1)^i(w/2 - t)\mathbf{v}\right) = 0.$$

Let $\mathbf{n}_i = \nabla f(\mathbf{c}_i)/\|\nabla f(\mathbf{c}_i)\|_2$ be 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 friction coefficient $\hat{\gamma}$, we discretize the friction cone at $\mathbf{c}_i$ [36] into a set of $l$ facets with vertices $\mathcal{F}_i = \left\{\mathbf{n}_i + \hat{\gamma}\cos\left(\frac{2\pi j}{l}\right)\mathbf{t}_{i,1} + \hat{\gamma}\sin\left(\frac{2\pi j}{l}\right)\mathbf{t}_{i,2} \big| j = 1, ..., l\right\}$. Each force $\mathbf{f}_{i,j} \in \mathcal{F}_i$ 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 $\mathbf{c}_i$. Under the soft contact model, each contact $\mathbf{c}_i$ exerts an additional wrench $\mathbf{w}_{i,l+1} = (\mathbf{0}, \mathbf{n}_i)$ [48]. Thus the set of all contact

wrenches that can be applied by a grasp $\mathbf{g}$ under the model is $\mathcal{W} = \{\mathbf{w}_{i,j} = (\mathbf{f}_{i,j}, \tau_{i,j})\big| i = 1, 2 \text{ and } j = 1, ..., l+1\}$.

### D. Quality Metric

In this work we use the probability of force closure ($P_F$), or the ability to resist external force and torques in arbitrary directions [30], as the quality metric. $P_F$ has shown promise in physical experiments [23], [44] and is relatively inexpensive to evaluate, allowing us to better study the effects of large amounts of data.

Let $F \in \{0, 1\}$ denote the occurrence of force closure. For a grasp $\mathbf{g}$ on object $\mathcal{O}$, $P_F(\mathbf{g}, \mathcal{O}) = \mathbb{P}(F = 1 \mid \mathbf{g}, \mathcal{O}, \xi, \nu, \gamma)$. To compute force closure for a grasp $\mathbf{g} \in \mathcal{G}$ on object $\mathcal{O} \in \mathcal{H}$ given samples of object pose $\hat{\xi}$, gripper pose $\hat{\nu}$, and friction coefficient $\hat{\gamma}$, we first compute the set of possible contact wrenches $\mathcal{W}$. Then $F = 1$ if $\mathbf{0} \in Conv(\mathcal{W})$, where $Conv(\cdot)$ denotes the convex hull [44].

### E. Objective

We are interested in finding a grasp $\mathbf{g}^*$ that maximizes $P_F(\mathbf{g})$ [23], [26], [30], [44] over a budgeted maximum number of samples $T$. To perform this as quickly as possible we maximize over the sum of the $P_F$ for all sampled grasps [26], [41]. Since the maximization over the continuous space $\mathcal{G}$ is computationally expensive, past work has solved this objective by evaluating a discrete set of $K$ candidate grasps $\Gamma = \{\mathbf{g}_1, ..., \mathbf{g}_K\}$ with Monte-Carlo integration [20], [44] or Multi-Armed Bandits (MAB) [26]. In this work, we extend the 2D MAB model of [26] to leverage similarities between grasps and prior 3D objects in Dex-Net to reduce the number of samples [16], [34].

## IV. DEXTERITY NETWORK

The Dexterity Network (Dex-Net) 1.0 dataset is a growing set that currently includes over 10,000 unique 3D object models annotated with 2.5 million parallel-jaw grasps.

### A. Data

Dex-Net 1.0 contains 13,252 3D mesh models: 8,987 from the SHREC 2014 challenge dataset [29], 2,539 from ModelNet40 [46], 1,371 from 3DNet [45], 129 from the KIT object database* [19], 120 from BigBIRD* [39], 80 from the Yale-CMU-Berkeley dataset* [7], and 26 from the Amazon Picking Challenge* scans (* indicates laser-scanner data). We preprocess each mesh by removing unreferenced vertices, computing a reference frame with Principal Component Analysis (PCA) on the mesh vertices, setting the mesh center of mass $\mathbf{z}$ to the center of the mesh bounding box, and rescaling the synthetic meshes to fit the smallest dimension of the bounding box within $w = 0.1m$. To resolve orientation ambiguity in the reference frame, we orient the positive $z$-axis toward the side of the $xy$ plane with more vertices. We also convert each mesh to an SDF using SDFGen [3].

## B. Grasp Sampling

Each 3D object $\mathcal{O}_i$ in Dex-Net is labelled with up to 250 parallel-jaw grasps and their $P_F$. We generate $N_g$ grasps for each object using a modification of the 2D algorithm presented in Smith et al. [40] to concentrate samples on grasps that are antipodal [30]. Let $w$ be the maximal opening of the gripper, $\hat{\gamma}$ be a sampled friction coefficient, and $\mathcal{S}$ be the set of points on the object surface for an SDF $f$ as described in Section III-A. To sample a single grasp, we first generate a contact point $\mathbf{c}_1$ by sampling uniformly from $\mathcal{S}$. Next we sample a direction $\mathbf{v} \in \mathbb{S}^2$ uniformly at random from the friction cone and compute $\mathbf{c}_2 = \mathbf{c}_1 + (w - t_2^*)\mathbf{v}$ and $\mathbf{x} = 0.5(\mathbf{c}_1 + \mathbf{c}_2)$, where $t_2^*$ is defined as in Section III-C. This yields a grasp $\mathbf{g}_{i,k} = (\mathbf{x}, \mathbf{v})$. We add $\mathbf{g}_{i,k}$ to the candidate set if the contacts are antipodal [30], or $\mathbf{v}^T \mathbf{n}_1 \leqslant \cos(\arctan(\hat{\gamma}))$ and $\mathbf{v}^T \mathbf{n}_2 \leqslant \cos(\arctan(\hat{\gamma}))$.

We evaluate $P_F(\mathbf{g})$ using Monte-Carlo integration [20] by sampling the object pose, gripper pose, and friction random variables $N_s$ times and recording $S_{i,k}$, the number of samples for which grasp $\mathbf{g}_{i,k}$ was in force closure. The dataset of $N_o$ objects and $N_g$ candidate grasps per object is $\mathcal{D} = \{(S_{i,k}, \mathcal{Y}_{i,k}) | i = \{1, ..., N_o\}, k = \{1, ..., N_g\}\}$ where $\mathcal{Y}_{i,k} = (\mathbf{g}_{i,k}, \mathcal{O}_i) \in \mathcal{M}$ is a grasp-object pair in Dex-Net.

## C. Grasp Differential Heightmap Features

To measure grasp similarity, we embed each grasp $\mathbf{g} = (\mathbf{x}, \mathbf{v})$ on object $\mathcal{O}$ in Dex-Net in a feature space based on 2D projections of the local surface orientation at the contacts, inspired by grasp heightmaps [15], [18]. Let $\delta \in \mathbb{R}$ be the pixel resolution in meters, and let $r \in \mathbb{R}$ be a minimum projection distance. The heightmap $\mathbf{h}_i$ at contact $\mathbf{c}_i$ maps points $\mathbf{p} \in \mathbb{R}^3$ on the tangent plane $\mathbf{v}^T(\mathbf{p} - \mathbf{c}_1) = 0$ to the distance to the surface along $\mathbf{v}$. To compute the value at pixel $u, v$, we compute the location of the pixel on the plane $\mathbf{p}_i(u,v) = \mathbf{c}_i + \delta u \mathbf{t}_1 + \delta v \mathbf{t}_2$ and assign

$$\mathbf{h}_i(u,v) = \min_{t \geqslant -r} t \text{ such that } f\left(\mathbf{p}_i(u,v) + (-1)^i t \mathbf{v}\right) = 0$$

where $f$ is the SDF of object $\mathcal{O}$. We make $\mathbf{h}_i$ rotation-invariant by orienting its axes to align with the eigenvectors of the weighted covariance matrix of the 3D surface points that generate the heightmap as described in [38]. Fig. 3 illustrates local surface patches extracted by this procedure. Since force closure depends on object surface normals at contacts [36], we finally take the $x$- and $y$-image gradients of $\mathbf{h}_i$ to form differential heightmaps $\mathbf{d}_{i,x}$ and $\mathbf{d}_{i,y}$. The feature vector for each grasp-object pair in Dex-Net is $\eta(\mathbf{g}, \mathcal{O}) = (\mathbf{d}_{1,x}, \mathbf{d}_{1,y}, \mathbf{d}_{2,x}, \mathbf{d}_{2,y})$.

## V. DEEP LEARNING FOR OBJECT SIMILARITY

We use Multi-View Convolutional Neural Networks (MV-CNNs) [43] to efficiently index prior 3D object and grasp data from Dex-Net by embedding each object in a vector space where distance represents object similarity, as shown in Fig. 4. We first render every object on a white background in a total of $N_c = 50$ virtual camera views oriented toward the object center and spaced on a grid of angle increments $\delta_\theta = \frac{2\pi}{5}$ and
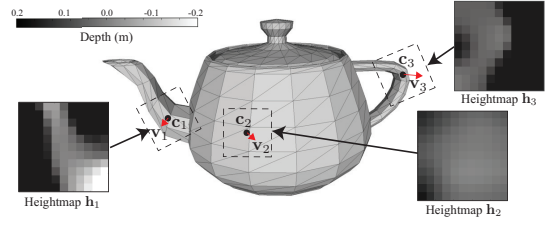


Fig. 3: Illustration of three local surface heightmaps extracted on a teapot. Each heightmap is "rendered" along the grasp axis $\mathbf{v}_i$ at contact $\mathbf{c}_i$ and oriented by the directions of maximum variation in the heightmap. We use gradients of the heightmaps for similiarity between grasps in Dex-Net.

$\delta_\varphi = \frac{2\pi}{5}$ on a viewing sphere with radii $r = R, 2R$, where $R$ is the maximum dimension of the object bounding box. Then we train a CNN with the architecture of AlexNet [24] to predict the 3D object class label for the rendered images on a training set of models. We initialize the weights of the network with the weights learned on ImageNet by Krizhevsky et al. [24] and optimize using Stochastic Gradient Descent (SGD). Next, we pass each of the $N_c$ views of each object through the optimized CNN and max-pool the output of the fc7 layer, the highest layer of the network before the class label prediction. Finally, we use Principal Component Analysis (PCA) to reduce the max-pooled output from 4,096 dimensions to 100 dimensions. This yields the representation $\psi(\mathcal{O}) \in \mathbb{R}^{100}$ for each object.

Given the MV-CNN object representation, we measure the dissimilarity between two objects $\mathcal{O}_i$ and $\mathcal{O}_j$ by the Euclidean distance $\|\psi(\mathcal{O}_i) - \psi(\mathcal{O}_j)\|_2$. For efficient lookups of similar objects, Dex-Net contains a KD-Tree nearest neighbor query structure with the feature vectors of all prior objects. In our implementation, we trained the MV-CNN using the Caffe library [17] on rendered images from a training set of approximately $6,000$ 3D models from the SHREC 2014 dataset [29], which has 171 unique categories, for 500,000 iterations of SGD. To validate the implementation, we tested on the SHREC 2014 challenge dataset and achieved a 1-NN accuracy of 86.7%, compared to 86.8% achieved by the winner of SHREC 2014 [29].

## VI. CORRELATED MULTI-ARMED BANDIT ALGORITHM

The Dex-Net 1.0 algorithm (see pseudocode below) optimizes $P_F$ over a set of candidate grasps on a new object $\mathcal{O}$ using Multi-Armed Bandits (MABs) with correlated rewards [16], [34] and priors computed from Dex-Net 1.0. We first generate a set of candidate grasps $\Gamma$ for object $\mathcal{O}$ using the antipodal grasp sampling described in Section IV-B and predict a prior belief distribution for each grasp using the Dex-Net database $\mathcal{D}$. Next, we run MAB by selecting a grasp using Thompson sampling [26], [33], sampling from the uncertainty random variables, determining force closure for the grasp on the sampled variables as described in Section III-D, and updating a belief distribution on the $P_F$ for each grasp. Finally, we rank the grasps in $\Gamma$ by the maximum lower confidence bound of the belief distribution, a conservative estimate of the $P_F$ of each grasp, and store the ranking in the database. We use Thompson sampling to study the scaling effects for a
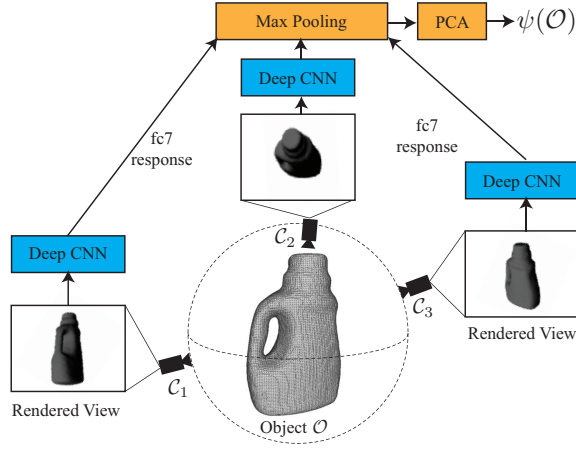
Fig. 4: Illustration of our Multi-View Convolutional Neural Network (MV-CNN) deep learning method for embedding 3D object models in a Euclidean vector space to compute global shape similarity. We pass a set of 50 virtually rendered camera viewpoints discretized around a sphere through a deep Convolutional Neural Network (CNN) with the AlexNet [24] architecture. Finally, we take the maximum fc7 response across each of the 50 views for each dimension and run PCA to reduce dimensionality.

fixed grasp selection method and plan to study other methods based on confidence bounds [25], [33] or Gittins indices [26] in future work.

### A. Belief Distribution Model

Let $\mathcal{O}$ denote the test object to label with the Dex-Net Algorithm, and let $\Gamma$ be the set of $N_g$ candidate grasps generated for $\mathcal{O}$. We define $F_j = F(\mathbf{g}_j) \in \{0, 1\}$ as force closure on an evaluation of any grasp $\mathbf{g}_j \in \Gamma$ from samples of object pose, gripper pose, and friction as described in Section III-D. Under the model, $F_j$ is a Bernoulli random variable with probability of success $\theta_j = P_F(\mathbf{g}_j)$. Since $\theta_j$ is unknown, the algorithm maintains a posterior Beta belief distribution on the Bernoulli parameter $\theta_j$ that is updated with every new observation of $F$, assigning increasingly high probability to the true $P_F$. The Beta distribution [16], [26] is specified by shape parameters $\alpha > 0$ and $\beta > 0$:

$$\text{Beta}(\alpha, \beta) = Z(\alpha, \beta)\theta_j^{\alpha-1}(1 - \theta_j)^{\beta-1}$$

where $Z(\alpha, \beta)$ is a normalization constant.

### B. Predicting Grasp Quality Using Prior Data

We use Continuous Correlated Beta Processes (CCBPs) [11], [32] to model correlations between the $P_F$ of grasps on different objects, which allows us to utilize prior grasp and object data from Dex-Net. CCBPs model correlations between Bernoulli random variables in a Beta-Bernoulli process, which exist when the variables depend on common latent factors. Two grasps on an object may have similar $P_F$ when they contact the object at similar locations, as evidenced by Lipschitz bounds on grasp wrench space metrics [36].

A CCBP estimates the shape parameters for a grasp-object pair $\mathcal{Y}_j = (\mathbf{g}_j, \mathcal{O}_i) \in \mathcal{M}$ 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.

We measure similarity using a set of feature maps $\varphi_m : \mathcal{M} \rightarrow \mathbb{R}^{d_m}$ for $m = 1, ..., 3$, where $d_m$ is the dimension of the feature space for each. The first feature map $\varphi_1(\mathcal{Y}) = (\mathbf{x}, \mathbf{v}, \|\rho_1\|_2, \|\rho_2\|_2)$ captures similiarity in the grasp parameters, where $\mathbf{x} \in \mathbb{R}^3$ is the grasp center, $\mathbf{v} \in \mathbb{S}^2$ is the grasp approach, and $\rho_i \in \mathbb{R}^3$ is the $i$-th moment arm. To capture local surface geometry, the second feature map $\varphi_2(\mathcal{Y}) = \eta(\mathbf{g}, \mathcal{O})$, where $\eta$ is the differential heightmap described in Section IV-C. To capture global shape information, our third feature map $\varphi_3(\mathcal{Y}) = \psi(\mathcal{O})$, where $\psi$ is our object similarity map described in Section V. Given the feature maps, we use the squared exponential kernel

$$k(\mathcal{Y}_p, \mathcal{Y}_q) = \exp\left(-\frac{1}{2}\sum_{m=1}^{3}\|\varphi_m(\mathcal{Y}_p) - \varphi_m(\mathcal{Y}_q)\|_{C_m}^2\right).$$

where $C_m \in \mathbb{R}^{d_m \times d_m}$ is the inverse bandwidth for $\varphi_m$ and $\|\mathbf{y}\|_{C_m} = \mathbf{y}^T C_m^T C_m \mathbf{y}$. The bandwidths are set by maximizing the log-likelihood [11] of the true $P_F$ under the CCBP on a set of training data.

We form a prior belief distribution for each candidate grasp in $\Gamma$ based on its similarity to all grasps and objects from the Dex-Net 1.0 database $\mathcal{D}$ as measured by the kernel [11]:

$$\alpha_{j,0} = \alpha_0 + \sum_{i=1}^{N_o}\sum_{k=1}^{N_g} k(\mathcal{Y}_j, \mathcal{Y}_{i,k})S_{i,k} \qquad \text{(VI.1)}$$

$$\beta_{j,0} = \beta_0 + \sum_{i=1}^{N_o}\sum_{k=1}^{N_g} k(\mathcal{Y}_j, \mathcal{Y}_{i,k})(N_s - S_{i,k}) \qquad \text{(VI.2)}$$

where $\alpha_0$ and $\beta_0$ are prior parameters for the Beta distribution [26] and $N_s$ is the number of times each grasp in $\mathcal{D}$ was sampled to evaluate $P_F$. In practice, we estimate the above sums using the $N_n$ nearest neighbors to $\mathcal{O}$ in the object similarity KD-Tree described in Section V. Upon observing $F_\ell$ for grasp $\mathbf{g}_\ell$ on iteration $t$, we update our belief for all other grasps on object $\mathcal{O}$ by [11]:

$$\alpha_{j,t} = \alpha_{j,t-1} + k(\mathcal{Y}_j, \mathcal{Y}_\ell)F_\ell \qquad \text{(VI.3)}$$

$$\beta_{j,t} = \beta_{j,t-1} + k(\mathcal{Y}_j, \mathcal{Y}_\ell)(1 - F_\ell). \qquad \text{(VI.4)}$$

## VII. EXPERIMENTS

We evaluate the convergence rate of the Dex-Net 1.0 algorithm for varying sizes of prior data used from Dex-Net and explore the sensitivity of the convergence rate to object shape, the similarity kernel bandwidths, and uncertainty. We created two training sets of 1,000, and 10,000 objects by uniformly sampling objects from Dex-Net. We uniformly sampled a set of 300 validation objects for selecting algorithm hyperparameters and selected a set of 45 test objects from the remaining objects. We ran the algorithm with $N_n = 10$ nearest neighbors, $\alpha_0 = \beta_0 = 1.0$ [26], and a lower confidence bound containing $p = 75\%$ of the belief distribution. We used isotropic Gaussian uncertainty with object and gripper translation variance $\sigma_t = 0.005$, object

**Input:** Object $\mathcal{O}$, Number of Candidate Grasps $N_g$, Number of Nearest Neighbors $N_n$, Dex-Net 1.0 Database $\mathcal{D}$, Features maps $\psi$ and $\eta$, Maximum Iterations $T$, Prior beta shape $\alpha_0$, $\beta_0$, Lower Bound Confidence $p$, Random Variables $\nu$, $\xi$, and $\gamma$

**Result**: Estimate of the grasp with highest $P_F$, $\hat{\mathbf{g}}^*$

```
    // Generate candidate grasps and priors
2   Γ = AntipodalGraspSample(O, N_g) ;
3   A_0 = ∅, B_0 = ∅;
4   for g_k ∈ Γ do
        // Equations VI.1 and  VI.2
5       α_{k,0}, β_{k,0} = ComputePriors(O, g_k, D, N_n, ψ);
6       A_0 = A_0 ∪ {α_{k,0}}, B_0 = B_0 ∪ {β_{k,0}};
7   end
    // Run MAB to Evaluate Grasps
8   for t = 1, .., T do
9       j = ThompsonSample(A_{t-1}, B_{t-1});
10      ν̂, ξ̂, γ̂ = SampleRandomVariables(ν, ξ, γ);
11      F_j = EvaluateForceClosure(g_j, O, ν̂, ξ̂, γ̂);
        // Equations VI.3 and  VI.4
12      A_t, B_t = UpdateBeta(j, F_j, Γ);
13      g_t^* = MaxLowerConfidence(A_t, B_t, p);
14  end
15  return g_T^*;
```

**Dex-Net 1.0 Algorithm**: Robust Grasp Planning Using Multi-Armed Bandits with Correlated Rewards
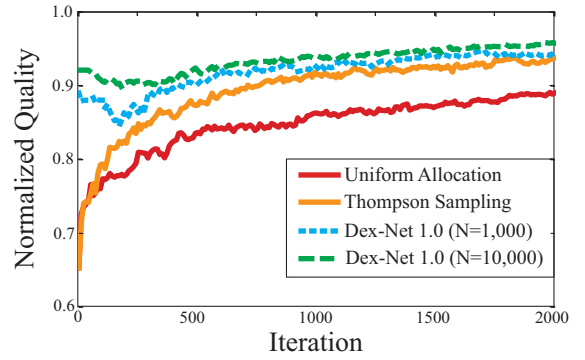


Fig. 5: Average normalized grasp quality versus iteration over 45 test objects and 25 trials per object for the Dex-Net1.0 algorithm with 1,000 and 10,000 prior 3D objects from Dex-Net. We measure quality by the $P_F$ for the best grasp predicted by the algorithm on each iteration and compare with Thompson sampling without priors and uniform allocation. The algorithm converges faster with 10,000 models, never dropping below approximately 90% of the grasp with highest $P_F$ from a set of 250 candidate grasps.

and gripper rotation variance $\sigma_r = 0.1$, and friction variance $\sigma_\gamma = 0.1$. For each experiment we compare the Dex-Net algorithm to Thompson sampling without priors (TS) and uniform allocation (UA) [26].

The inverse kernel bandwidths were selected by maximizing the log-likelihood of the true $P_F$ under the CCBP model [11] on the validation set using a grid search over hyperparameters. The inverse bandwidths of the similarity kernel were $C_g = diag(0, 0, 175, 175)$ for the grasp parameter features, an isotropic Gaussian mask $C_d$ with mean $\mu_d = 500.0$ and $\sigma_d = 1.75$ for the differential heightmap features, and $C_s = 0.001 * I$ for the shape similarity features.

To scale the experiments, we developed a Cloud-based system on top of Google Cloud Platform. We used Google Compute Engine (GCE) to distribute trials of MAB algorithms across objects and Google Cloud Storage to store Dex-Net. The system launched up to 1,500 GCE virtual instances at once for experiments, reducing the runtime by three orders of magnitude. Each virtual instance ran Ubuntu 12.04 on a single core with 3.75 GB of RAM.

### A. Scaling of Average Convergence Rate

To examine the effects of orders of magnitude of prior data on convergence to a grasp with high $P_F$, we ran the Dex-Net 1.0 algorithm on the test objects with priors computed from 1,000 and 10,000 objects from Dex-Net. Fig. 5 shows the normalized $P_F$ (the ratio of the $P_F$ for the sampled grasp to the highest $P_F$ of the 250 candidate grasps) versus iteration averaged over 25 trials for each of the test objects over 2,000 iterations. The average runtime per iteration was 16 ms for UA, 17 ms for TS, and 22 ms for Dex-Net 1.0. The algorithm with 10,000 objects takes approximately $2\times$ fewer iterations to reach the maximum normalized $P_F$ value reached by TS. Furthermore, the 10,000 object curve does not fall below

approximately 90% of the best grasp in the set across all iterations, suggesting that a grasp with high $P_F$ is found using prior data alone.

### B. Sensitivity to Object Shape

To understand the behavior of the Dex-Net algorithm on individual 3D objects, we examined the convergence rate with a 3D model of a drill and spray bottle from the test set, both uncommon object categories in Dex-Net. Fig. 1 and Fig. 6 show the normalized $P_F$ versus iteration averaged over 25 trials for 2,000 iterations on the spray bottle and drill, respectively. We see that the spray bottle converges very quickly when using a prior dataset of 10,000 objects, finding the optimal grasp in the set in about 1,500 iterations. This convergence may be explained by the two similar spray bottles retrieved by the MV-CNN from the 10,000 object dataset. Fig. 7 illustrates the grasps predicted to have the highest $P_F$ on the spray bottle by the different algorithms after 100 iterations. On the other hand, performance on the drill does not increase using either 1,000 or 10,000 objects, as the closest model in all of Dex-Net according to the similarity metric is a phone.

### C. Sensitivity to Similarity and Uncertainty

We also studied the sensitivity of the Dex-Net algorithm to the kernel bandwidth hyperparameters described in Section VI-B and the levels of pose and friction uncertainty for the test object. We varied the inverse bandwidths of the kernel for the grasp parameters and differential heightmaps gradients to the lower values $C_g = diag(0, 0, 15, 15)$, $\mu_d = 350.0$, and $\sigma_d = 3.0$ as well as the higher values $C_g = diag(0, 0, 300, 300)$, $\mu_d = 750.0$, and $\sigma_d = 1.75$. We also tested low uncertainty with variances $(\sigma_t, \sigma_r, \sigma_\gamma) = (0.0025, 0.05, 0.05)$ and high uncertainty with variances $(\sigma_t, \sigma_r, \sigma_\gamma) = (0.01, 0.2, 0.2)$. Fig. 8 shows the normalized $P_F$ versus iteration averaged over 25 trials for 2,000 iterations on the 45 test objects. The results suggest that conservative setting of similiarity kernel bandwidth is important for convergence and that the algorithm is not sensitive to uncertainty levels.
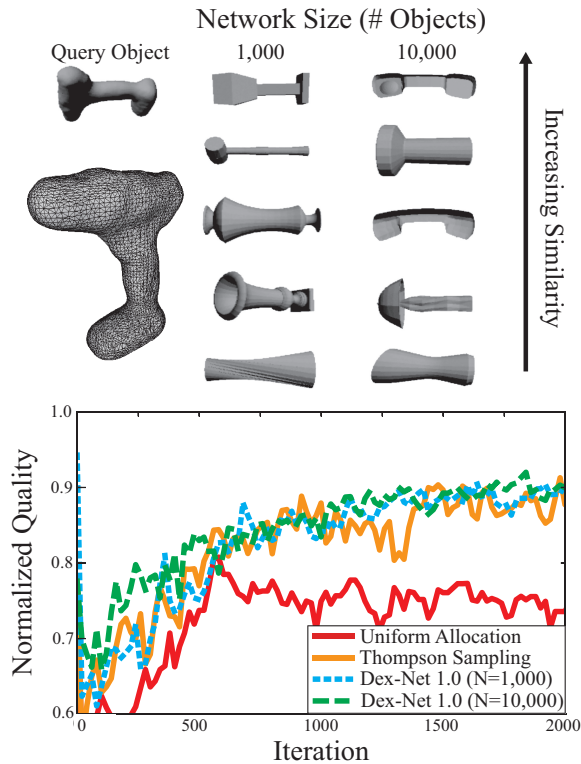
Fig. 6: Failure object for the Dex-Net 1.0 algorithm. (Top) The drill, which is relatively rare in the dataset, has no geometrically similar neighbors even with 10,000 objects. (Bottom) Plotted is the average normalized grasp quality versus iteration over 25 trials for the Dex-Net 1.0 algorithm with 1,000 and 10,000 prior 3D objects. The lack of similar objects leads to no significant performance increase over Thompson sampling without priors.
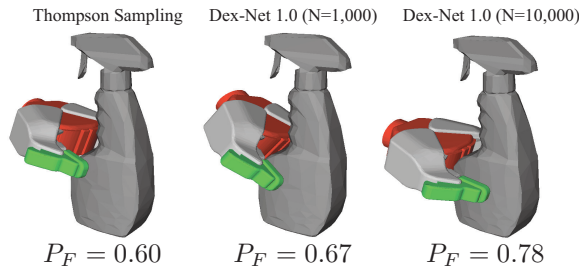


Fig. 7: Illustration of the grasps predicted to have the highest $P_F$ after only 100 iterations by Thompson sampling without priors and the Dex-Net 1.0 algorithm with 1,000 and 10,000 prior objects. Thompson sampling without priors chooses a grasp near the edge of the object, while the Dex-Net algorithm selects grasps closer to the object center-of-mass.

## VIII. DISCUSSION AND FUTURE WORK

We presented Dexterity Network 1.0 (Dex-Net), a new dataset and associated algorithm to study the scaling effects of Big Data and Cloud Computation on robust grasp planning. The algorithm uses a Multi-Armed Bandit model with correlated rewards to leverage prior grasps and 3D object models and Multi-View Convolutional Neural Networks (MV-CNNs), a new deep learning method for 3D object classification, as a similarity metric between objects. In experiments, the Google Cloud Platform allowed Dex-Net 1.0 to simultaneously run up to 1,500 virtual machines, reducing experiment runtime by three orders of magnitude. Experiments suggest that prior data can speed robust grasp planning by a factor of 2 and that
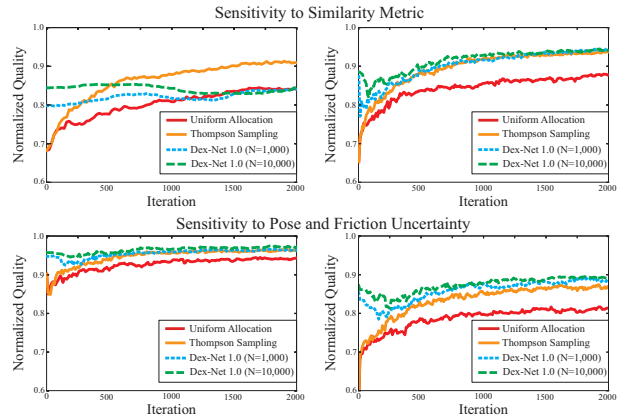


Fig. 8: Sensitivity to similiarity kernel (top) and pose and friction uncertainty (bottom) for the normalized grasp quality versus iteration averaged over 25 trials per object for the Dex-Net algorithm with 1,000 and 10,000 prior 3D objects. (Top-left) Using a higher inverse bandwidth causes the algorithm to measure false similarities between grasps, leading to performance on par with uniform allocation. (Top-right) A lower inverse bandwith decreases the convergence rate, but on average the Dex-Net Algorithm still selects a grasp within approximately 85% of the grasp with highest $P_F$ for all iterations. (Bottom-left) Lower uncertainty increases the quality for all methods and (bottom-right) higher uncertainty decreases the quality for all methods, and the Dex-Net algorithm with 10,000 prior objects still converges approximately $2\times$ faster than Thompson sampling without priors.

average grasp quality increases with the number of similar objects in the dataset. We reported on sensitivity to varying similarity metrics and pose and friction uncertainty levels.

In future work, we will develop metrics to pre-compute grasps that adequately "cover" each object from a variety of accessibility conditions (depending on pose and occlusions). We will also explore how Deep Learning [24] can be used in other parts of a grasp planning pipeline, for example to recognize object pose and shape from images [1], to learn grasp and object features robust to shape variation using prior evaluations from bandit algorithms, and perhaps even to determine motor torques based on images and precomputed grasps [28]. We also hope to release subsets of Dex-Net 1.0 with an open-source API to explore robust grasping as a service (RGaaS).

## REFERENCES

[1] M. Aubry and B. Russell, "Understanding deep features with computer-generated imagery," *arXiv preprint arXiv:1506.01151*, 2015.

[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] C. Batty, "Sdfgen," https://github.com/christopherbatty/SDFGen.

[4] J. Bohg, A. Morales, T. Asfour, and D. Kragic, "Data-driven grasp synthesis–a survey," *Robotics, IEEE Transactions on*, vol. 30, no. 2, pp. 289–309, 2014.

[5] A. M. Bronstein, M. M. Bronstein, L. J. Guibas, and M. Ovsjanikov, "Shape google: Geometric words and expressions for invariant shape retrieval," *ACM Transactions on Graphics (TOG)*, vol. 30, no. 1, p. 1, 2011.

[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] D.-Y. Chen, X.-P. Tian, Y.-T. Shen, and M. Ouhyoung, "On visual similarity based 3d model retrieval," in *Computer graphics forum*, vol. 22, no. 3. Wiley Online Library, 2003, pp. 223–232.

[9] 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.

[10] R. Detry, D. Kraft, O. Kroemer, L. Bodenhagen, J. Peters, N. Krüger, and J. Piater, "Learning grasp affordance densities," *Paladyn, Journal of Behavioral Robotics*, vol. 2, no. 1, pp. 1–17, 2011.

[11] 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.

[12] C. Goldfeder and P. K. Allen, "Data-driven grasping," *Autonomous Robots*, vol. 31, no. 1, pp. 1–20, 2011.

[13] 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.

[14] 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.

[15] 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.

[16] M. W. Hoffman, B. Shahriari, and N. de Freitas, "Exploiting correlation and budget constraints in bayesian multi-armed bandit optimization," *arXiv preprint arXiv:1303.6746*, 2013.

[17] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.

[18] D. Kappler, J. Bohg, and S. Schaal, "Leveraging big data for grasp planning," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2015.

[19] 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.

[20] B. Kehoe, D. Berenson, and K. Goldberg, "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.

[21] 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.

[22] B. Kehoe, S. Patil, P. Abbeel, and K. Goldberg, "A survey of research on cloud robotics and automation," *Automation Science and Engineering, IEEE Transactions on*, vol. 12, no. 2, pp. 398–409, 2015.

[23] 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.

[24] 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.

[25] 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.

[26] 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.

[27] 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.

[28] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *arXiv preprint arXiv:1504.00702*, 2015.

[29] 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.

[30] 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.

[31] D. Maturana and S. Scherer, "Voxnet: A 3d convolutional neural network for real-time object recognition," in *Proc. IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, 2015.

[32] 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.

[33] J. Oberlin and S. Tellex, "Autonomously acquiring instance-based object models from experience," 2013.

[34] S. Pandey, D. Chakrabarti, and D. Agarwal, "Multi-armed bandit problems with dependent arms," in *Proceedings of the 24th international conference on Machine learning*. ACM, 2007, pp. 721–728.

[35] F. T. Pokorny, K. Hang, and D. Kragic, "Grasp moduli spaces." in *Robotics: Science and Systems*, 2013.

[36] F. T. Pokorny and D. Kragic, "Classical grasp quality evaluation: New theory and algorithms," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2013.

[37] 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.

[38] S. Salti, F. Tombari, and L. Di Stefano, "Shot: Unique signatures of histograms for surface and texture description," *Computer Vision and Image Understanding*, vol. 125, pp. 251–264, 2014.

[39] A. Singh, J. Sha, K. S. Narayan, T. Achim, and P. Abbeel, "Bigbird: A large-scale 3d database of object instances," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2014.

[40] G. Smith, E. Lee, K. Goldberg, K. Bohringer, and J. Craig, "Computing parallel-jaw grips," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 1999.

[41] N. Srinivas, A. Krause, S. Kakade, and M. Seeger, "Gaussian process optimization in the bandit setting: No regret and experimental design," in *Proc. International Conference on Machine Learning (ICML)*, 2010.

[42] 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.

[43] 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.

[44] 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.

[45] W. Wohlkinger, A. Aldoma, R. B. Rusu, and M. Vincze, "3dnet: Large-scale object class recognition from cad models," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2012.

[46] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shape modeling," in *CVPR*, vol. 1, no. 2, 2015, p. 3.

[47] L. E. Zhang, M. Ciocarlie, and K. Hsiao, "Grasp evaluation with graspable feature matching," in *RSS Workshop on Mobile Manipulation: Learning to Manipulate*, 2011.

[48] 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.