

GraspNet: A Large-Scale Clustered and Densely Annotated Dataset for Object Grasping

Hao-Shu Fang, Chenxi Wang, Minghao Gou, Cewu Lu
Shanghai Jiao Tong University

fhaoshu@gmail.com, {wcx1997, gmh2015, lucewu}@sjtu.edu.cn

Abstract

Object grasping is critical for many applications, which is also a challenging computer vision problem. However, for clustered scene, current researches suffer from the problems of insufficient training data and the lacking of evaluation benchmarks. In this work, we contribute a large-scale grasp pose detection dataset with an unified evaluation system. Our dataset contains 87,040 RGBD image with over 370 million grasp poses. Meanwhile, our evaluation system directly reports whether a grasping is successful or not by analytic computation, which is able to evaluate any kind of grasp poses without exhausted labeling pose ground-truth. We conduct extensive experiments to show that our dataset and evaluation system can align well with real-world experiments. Our dataset, source code and models will be made publicly available.

1. Introduction

Object grasping is a fundamental problem and has many applications in industry, agriculture and service trade. The key of grasping is to detect the grasp pose given visual inputs (image or point cloud) and has drawn many attentions in computer vision community [8, 21].

Though important, there are currently two main hindrances to obtain further performance gains in this area. Firstly, the grasp pose has different representations including rectangle [23] and 6D pose [24] representation and are evaluated with different metrics [11, 10, 24] correspondingly. The difference in evaluation metrics makes it difficult to compare these methods directly in an unified manner, while evaluating with real robots would dramatically increase the evaluation cost. Secondly, it is difficult to obtain large-scale high quality training data [3]. Previous datasets annotated by human [11, 27, 5] are usually small in scale and only provide sparse annotations. While obtaining training data from the simulated environment [17, 7, 26] can generate large scale datasets, the visual domain gap be-

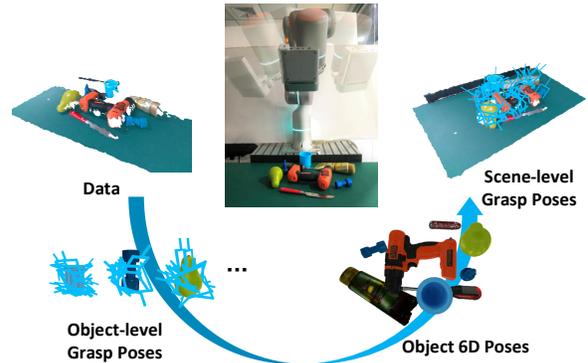


Figure 1. Our methodology for building the dataset. We collect data with real-world sensors and annotate grasp poses for every single object by analytic computation. Object 6D poses are manually annotated to project the grasp poses from object coordinate to the scene coordinate. Such methodology greatly reduces the labor of annotating grasp poses. Our dataset is both densely annotated and visually coherent with real world.

tween simulation and reality would inevitably degrade the performance of algorithms in real-world application.

To form a solid foundation for algorithms built upon, it is important for a dataset to i) provide data that aligns well with the visual perception from real world sensor, ii) be densely and accurately annotated with large-scale grasp pose ground-truth and iii) evaluate grasp poses with different representations in a unified manner. This is nontrivial, especially when it comes to the data annotation. Given an image or scene, it's hard for us to manually annotate endless grasp poses in continuous space. We circumvent this issue by exploring a new direction, that is, collecting data from the real world and annotating them by analytic computation in simulation, which leverages the advantages from both sides.

Specifically, inspired by previous literature [24], we propose a two-step pipeline to generate tremendous grasp poses for a scene. Thanks to our automatic annotation process, we built the first large-scale in-the-wild grasp pose dataset that can serve as a base for training and evaluating grasp pose

detection algorithms. Our dataset contains **87,040** RGB-D images taken from different viewpoints of over 170 clustered scenes.

For all 88 objects in our dataset, we provide accurate 3D mesh models. Each scene is densely annotated with object 6D pose and grasp pose, resulting in over **370 million** grasp poses, which is 3 orders of magnitude larger than previous datasets. Moreover, embedded with an online evaluation system, our benchmark is able to evaluate current mainstream grasping detection algorithms. Experiments also demonstrate that our benchmark can align well with real-world experiments. Fig 1 shows the methodology for building our dataset.

2. Related Work

In this section, we first review deep learning based grasping detection algorithms, followed by related datasets in this area.

Deep Learning Based Grasping Prediction Algorithms

For deep learning based grasping detection algorithms, they can be divided into three main categories. The most popular one is to detect a graspable rectangle based on RGB-D image input [11, 13, 23, 9, 21, 14, 17, 27, 1, 2, 5, 18, 19]. Lenz *et al.* [13] proposed a cascaded method with two networks that first prunes out unlikely grasps and then evaluates the remaining grasps with a larger network. Redmon *et al.* [23] proposed different network structure that directly regresses the grasp poses in a single step manner, which is faster and more accurate. Mahler *et al.* [17] proposed a grasp quality CNN to predict the robustness score of grasping candidates. Zhang [27] and Chu [5] extended it to multi-object scenarios. The grasp poses generated by these methods are constrained in 2D plane which limits the degree of freedom of grasp poses. With the rapid development in monocular object 6D pose estimation [12, 25, 28], some researchers [6] predicted 6D pose of the objects and project predefined grasp pose on the scene. Such methods have no limitation of grasping orientation, but require a prior knowledge about the object shape. Recently, there is a new line of researches [24, 15, 19, 22] that proposed grasping candidates on partial observed point cloud and output a classification score for each candidate using 3D CNN. Such methods have no limitation and require no prior knowledge about the objects. Currently, these methods are evaluated in their own metrics and hard to compare to others.

Grasping Dataset Cornell grasping dataset [11] first proposed rectangle representation for grasping detection in images. Single object RGB-D images are provided with rectangle grasp poses. [5, 27] built datasets with the same protocol but extend to multi-object scenarios. These grasp

poses are annotated by human. [21, 14] collected annotations with real robot experiments. These data labeling methods are time consuming and require strong hardware support. To avoid such problem, some recent works explored using simulated environment [17, 7, 26, 19] to annotate grasp poses. They can generate a much larger scale dataset but the domain gap of visual perception is always a hindrance. Beyond rectangle based annotation, GraspSeg [2] provided pixel-wise annotations for grasp-affordance segmentation and object segmentation. For 6D pose estimation, [25] contributed a dataset with 21 objects and 70 scenes. These datasets mainly focused on a subarea of grasp pose detection. In this work, we aim to build a dataset that is much larger in scale and diversity and covers main aspects of object grasping detection.

3. GraspNet Dataset

We next describe the main features of our dataset and how we build it.

3.1. Overview

Previous grasping dataset either focuses on isolated object [11, 17, 7, 26] or only labels one grasp per scene [21, 14]. Few datasets consider multi-object-multi-grasp setting and are small in scale [27, 5] due to the labor of annotation. Moreover, most of the datasets adopt the rectangle based representation [11] of grasp pose, which constrains the space for placing the gripper. To overcome these issues, we propose a large-scale dataset in clustered scenario with dense and rich annotations for grasp pose prediction named *GraspNet*. The GraspNet dataset contains 88 daily objects with high quality 3D mesh models. The images are collected from 170 clustered scenes, each contributes 512 RGB-D images captured by two different cameras, resulting 87,040 images in total. For each image, we densely annotate 6-DoF grasp poses by analytic computation of force closure [20]. The grasp poses for each scene varies from 1,500,000 to 2,500,000, and in total our dataset contains over 370 million grasp poses. Besides, we also provide accurate object 6D pose annotations, rectangle based grasp poses, object masks and bounding boxes. Each frame is also associated with a camera pose, thus multi-view point cloud can be easily fused. Fig 2 illustrates the key components of our dataset.

3.2. Data Collection

We select 32 objects that are suitable for grasping from the YCB dataset [4], 13 adversarial objects from DexNet 2.0 [17] and collect 43 objects of our own to construct our object sets. The objects vary in shape, texture, size and material. To collect data of clustered scene, we attach the cameras to a robot arm since it can repeat the trajectory precisely

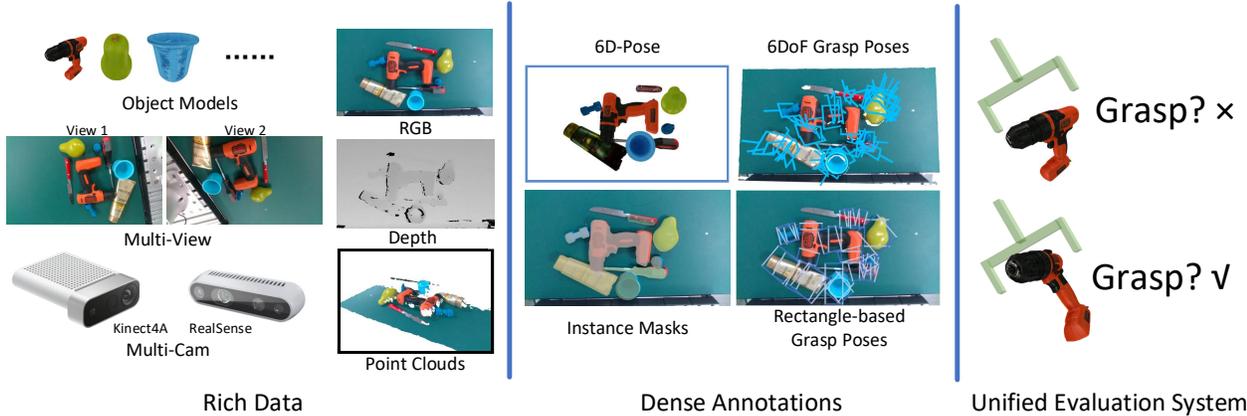


Figure 2. The key components of our dataset. RGB-D images are taken using both RealSense camera and Kinect camera from different views. The 6D pose of each object, the grasp poses, the rectangle grasp poses and the instance masks are annotated. A unified evaluation system is also provided.

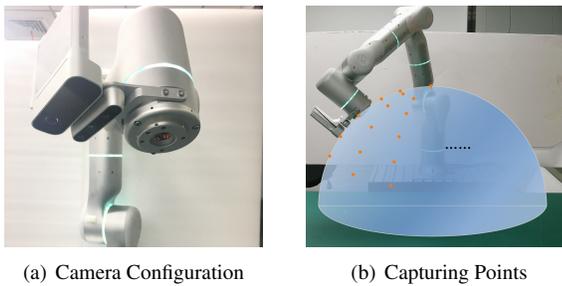


Figure 3. The setting of data collection. (a) Both the RealSense and Kinect camera are fixed on the end link of a robot arm. (b) 256 data collection points are sampled from a quarter sphere.

and help automatizing the collecting process. Hand-eye calibration is conducted before data collection to obtain accurate camera poses. Considering different quality of depth image will inevitably affect the algorithms, we adopt two popular RGB-D cameras, Intel RealSense 435 and Kinect 4 Azure, to simultaneously capture the scene and provide rich data. For each cluster scene, we randomly pick around 10 objects from our whole set and place them in a clustered manner. The robot arm then moves along a fixed trajectory that covers 256 distinct viewpoints on a quarter sphere. A synchronized image pair from both RGB-D cameras as well as their camera poses will be saved. Fig 3 shows the setting for our data collection.

3.3. Data Annotation

6D Pose Annotation With 87,040 images in total, it would be labor consuming to annotate 6D poses for each frame. Thanks to the camera poses recorded, we only need to annotate 6D poses for the first frame of each scene. The 6D poses will then be propagated to the remaining frames by:

$$\mathbf{P}_i^j = \mathbf{cam}_i^{-1} \mathbf{cam}_0 \mathbf{P}_0^j, \quad (1)$$

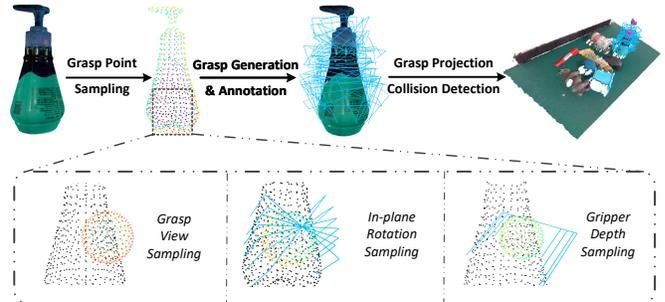


Figure 4. Grasp pose annotation pipeline. The grasp point is firstly sampled from point cloud. Then the grasp view, the in-place rotation angle and the gripper depth are sampled. Grasps with high scores are adopted for each object. Finally, the grasps are projected on the scene using the 6D pose of each object. Collision detection is also conducted to avoid the collision between grasps and background or other object.

where \mathbf{P}_i^j is the 6D pose of object j at frame i and \mathbf{cam}_i is the camera pose of frame i . Followed by ICP refinement and human checking, we can annotate 6D poses accurately in an efficient manner. Object masks and bounding boxes are also obtained by projecting objects onto the images using 6D poses.

Grasp Pose Annotation Different from labels in common vision tasks, grasp poses distribute in a large and continuous search space, which brings infinite annotations. Annotating each scene manually would be dramatically labor expensive. Considering all the objects are known, we propose a two stage automated pipeline for grasp pose annotation, which is illustrated in Fig. 4.

First, grasp poses are sampled and annotated for each single object. To achieve that, high quality mesh models are downsampled such that the sampled points (called grasp points) are uniformly distributed in voxel space. For each

Dataset	Grasps / scene	Objects / scene	Grasp label	6D pose	Total objects	Total grasps	Total images	Modality	Data source
Cornell [11]	~8	1	Rect.	No	240	8019	1035	RGB-D	1 Cam.
Pinto <i>et al.</i> [21]	1	-	Rect.	No	150	50K	50K	RGB-D	1 Cam.
Levine <i>et al.</i> [14]	1	-	Rect.	No	-	800K	800K	RGB-D	1 Cam.
Mahler <i>et al.</i> [17]	1	1	Rect.	No	1,500	6.7M	6.7M	Depth	Sim.
Jacquard [7]	~20	1	Rect.	No	11K	1.1M	54K	RGB-D	Sim.
Zhang <i>et al.</i> [27]	~20	~3	Rect.	No	-	100K	4683	RGB	1 Cam.
Multi-Object [5]	~30	~4	Rect.	No	-	2904	96	RGB-D	1 Cam.
VR-Grasping-101* [26]	100	1	6-DOF	Yes	101	4.8M	10K	RGB-D	Sim.
YCB-Video [25]	None	~5	None	Yes	21	None	134K	RGB-D	1 Cam.
GraspNet (ours)	~2.2M	~10	6-DOF	Yes	88	370M	87K	RGB-D	2 Cams.

Table 1. Summary of the properties of publicly available grasp datasets. “Rect.,” “Cam.” and “Sim.” are short for Rectangle, Camera and Simulation respectively. “-” denotes the number is unknown. “*” denotes the dataset is temporarily unavailable according to the author.

grasp point, we sample V views uniformly distributed in a spherical space. Grasp candidates are searched in a two dimensional grid $D \times A$, where D is the set of gripper depths and A is the set of in-plane rotation angles. Gripper width is determined accordingly such that no empty grasp or collision occurs. Each grasp candidate will be assigned a confidence score based on the mesh model.

We adopt an analytic computation method to grade each grasp. The force-closure metric [20, 24] has been proved effective in grasp evaluation: given a grasp pose, the associated object and a friction coefficient μ , force-closure metric outputs a binary label indicating whether the grasp is antipodal under that coefficient. The result is computed based on physical rules, which is robust. Here we adopt an improved metric described in [15]. With $\Delta\mu = 0.1$ as interval, we increase μ gradually from 0.1 to 1 step by step until the grasp is antipodal. The grasp with lower friction coefficient μ has more probability of success. Thus we define our score s as:

$$s = 1.1 - \mu, \quad (2)$$

such that s lies in $(0, 1]$.

Second, for each scene, we project these grasps to the corresponding objects based on the annotated 6D object poses:

$$\begin{aligned} \mathbf{P}^i &= \mathbf{cam}_0 \mathbf{P}_0^i, \\ \mathbb{G}_{(w)}^i &= \mathbf{P}^i \cdot \mathbb{G}_{(o)}^i, \end{aligned} \quad (3)$$

where \mathbf{P}^i is the 6D pose of the i -th object in the world frame, $\mathbb{G}_{(o)}^i$ is a set of grasp poses in the object frame and $\mathbb{G}_{(w)}^i$ contains the corresponding poses in the world frame. Besides, collision check is performed to avoid invalid grasps. After these two steps we can generate densely distributed grasp set $\mathbb{G}_{(w)}$ for each scene. According to statistics, the ratio of positive and negative labels in our dataset is around 1:2. We conduct real world experiment in Sec. 4 using our robot arm and verify that our generated grasp poses can align well with real world grasping.

3.4. Evaluation

Dataset Split For our 170 scenes, we use 100 for training and 70 for testing. Specifically, we further divide our test sets into 3 categories: 30 scenes with seen objects, 30 with unseen but similar objects and 10 for novel objects. We hope that such setting can better evaluate the generalization ability of different methods.

New Metrics To evaluate the prediction performance of grasp pose, previous methods adopt the rectangle metric that consider a grasp as correct if: i) the rotation error is less than 30° and ii) the rectangle IOU is larger than 0.25.

There are several drawbacks of such metric. Firstly, it can only evaluate rectangle representation of grasp pose. Secondly, the error tolerance is set rather high since the groundtruth annotations are not exhaustive. It might overestimate the performance of grasping algorithm. Currently, the Cornell dataset [11] has achieved over 99% accuracy. In this work, we adopt an online evaluation algorithm to evaluate the grasp accuracy.

We first illustrate how we classify whether a single grasp pose is true positive. For each predicted grasp pose $\hat{\mathbf{P}}_i$, we associate it with the target object by checking the point cloud inside the gripper. Then, similar to the process of generating grasp annotation, we can get a binary label for each grasp pose by force-closure metric, given different μ .

For clustered scene, grasp pose prediction algorithms are expected to predict multiple grasps. Since for grasping, we usually conduct execution after the prediction, the percentage of true positive is more important. Thus, we adopt *Precision@k* as our evaluation metric. *AP* denotes the average *Precision@k* for k ranges from 1 to 50. Similar to COCO [16], we report *AP* at different μ . Specifically, we denote **AP** for *AP* from $\mu = 0.1$ to $\mu = 0.5$, with $\Delta\mu = 0.1$ as interval.

To avoid dominated by similar grasp poses or grasp poses from single object, we run a pose-NMS before eval-

uation. For two grasps \mathbf{G}_1 and \mathbf{G}_2 , we define grasp pose distance $D(\mathbf{G}_1, \mathbf{G}_2)$ as a tuple:

$$D(\mathbf{G}_1, \mathbf{G}_2) = (d_t(\mathbf{G}_1, \mathbf{G}_2), d_\alpha(\mathbf{G}_1, \mathbf{G}_2)), \quad (4)$$

where $d_t(\mathbf{G}_1, \mathbf{G}_2)$ and $d_\alpha(\mathbf{G}_1, \mathbf{G}_2)$ denote translation distance and rotation distance of two grasps respectively. Let a grasp pose \mathbf{G} be denoted by a translation vector \mathbf{t} and a rotation matrix \mathbf{R} , then $d_t(\cdot)$ and $d_\alpha(\cdot)$ is defined as:

$$\begin{aligned} d_t(\mathbf{G}_1, \mathbf{G}_2) &= \|\mathbf{t}_1 - \mathbf{t}_2\|, \\ d_\alpha(\mathbf{G}_1, \mathbf{G}_2) &= \arccos \frac{1}{2}(\text{tr}(R_1 \cdot R_2^T) - 1), \end{aligned} \quad (5)$$

where $\text{tr}(\mathbf{M})$ denotes the trace of matrix \mathbf{M} .

Since translation and rotation are not in the same metric space, we define the NMS threshold as a tuple too. Let $TH = (th_d, th_\alpha)$, we say $D(\mathbf{G}_1, \mathbf{G}_2) < TH$ if and only if

$$d_t(\mathbf{G}_1, \mathbf{G}_2) < th_d, \quad d_\alpha(\mathbf{G}_1, \mathbf{G}_2) < th_\alpha. \quad (6)$$

Based on the tuple metric, two grasps are merged when their distance is lower than TH . Meanwhile, only the top K grasps from each object are considered according to confidence scores and other grasps are omitted. In evaluation, we set $th_d = 1$ cm, $th_\alpha = 5$ degree and $K = 10$.

3.5. Discussion

We compare our datasets with other publicly available grasp datasets. Table 1 summaries the main differences at several aspects. We can see that our dataset is much larger in scale and diversity. With our two-step annotation pipeline, we are able to collect real images with dense annotations, which leverages the advantages from both sides.

For grasp pose evaluation, due to the continuity in grasping space, there are in fact infinite feasible grasp poses. The previous method that pre-computed ground truth for evaluating grasping, no matter collected by human annotation [11] or simulation [7], cannot cover all feasible solution. In contrast, we do not pre-compute labels for the test set, but directly evaluate them by calculating the quality score using force closure metric [20]. Such evaluation method does not assume the representation of the grasp pose, thus is general in practice. Related APIs will be made publicly available to facilitate the research in this area.

4. Experiments

In this section, we conduct robotic experiments to demonstrate that our ground-truth annotations can align well with real-world grasping.

4.1. Ground-Truth Evaluation

To evaluate the quality of our generated grasp poses, we set up a real robotic experiment. Since we need to project

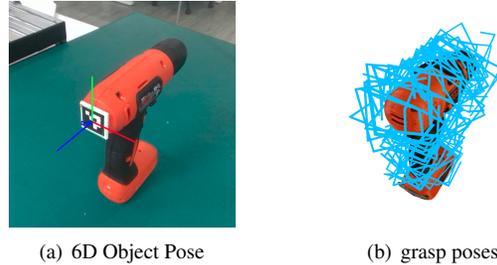


Figure 5. (a) Visualization of the detected 6D pose the ArUco marker. (b) Object 6D pose and grasp poses inferred from the marker pose.

Object	s=1	s=0.5	s=0.1	Object	s=1	s=0.5	s=0.1
Banana	98%	67%	21%	Apple	97%	65%	16%
Peeler	95%	59%	9%	Dragon	96%	60%	9%
Mug	96%	62%	12%	Camel	93%	67%	23%
Scissors	89%	61%	5%	Power Drill	96%	61%	14%
Lion	98%	68%	16%	Black Mouse	98%	64%	13%

Table 2. Summary of real world success rate of grasping given different grasp score.

grasp poses to the camera frame using objects' 6D poses, we paste ArUco code on the objects and only label their 6D poses once to avoid tedious annotation process. Fig 5 illustrates our grasp pose projection process.

We pick 10 objects from our object set and execute grasp poses that has different scores. For each setting we randomly choose 100 grasp poses. For robot arm we adopt a Flexiv Rizon arm and for camera we use the Intel RealSense 435. Table 2 summarizes the success rate of grasping. We can see that for grasp poses with high score, the success rate can achieve 0.96 in average. Meanwhile, the success rate is pretty low for grasp poses with $s = 0.1$. It indicates that our generated grasp poses are well aligned with real world grasping.

5. Conclusion

In this paper we built a large-scale dataset for clustered scene object grasping. Our dataset is orders of magnitude larger than previous grasping datasets and diverse in objects, scenes and data sources. It consists of images taken by real world sensor and has rich and dense annotations. Meanwhile, a unified evaluation system is also proposed to promote the development of this area. We demonstrated that our dataset and evaluation system align well with real world grasping. In the future, we will also extend our dataset to multi-finger gripper and vacuum-based end effectors. Our code and dataset will be released.

References

- [1] Umar Asif, Jianbin Tang, and Stefan Herrer. Ensembled: Improving grasp detection using an ensemble of convolutional neural networks. In *BMVC*, page 10, 2018. 2
- [2] Umar Asif, Jianbin Tang, and Stefan Herrer. Graspnet: An efficient convolutional neural network for real-time grasp detection for low-powered devices. In *IJCAI*, pages 4875–4882, 2018. 2
- [3] Shehan Caldera, Alexander Rassau, and Douglas Chai. Review of deep learning methods in robotic grasp detection. *Multimodal Technologies and Interaction*, 2(3):57, 2018. 1
- [4] Berk Calli, Arjun Singh, James Bruce, Aaron Walsman, Kurt Konolige, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M Dollar. Yale-cmu-berkeley dataset for robotic manipulation research. *The International Journal of Robotics Research*, 36(3):261–268, 2017. 2
- [5] Fu-Jen Chu, Ruinian Xu, and Patricio A Vela. Real-world multiobject, multigrasp detection. *IEEE Robotics and Automation Letters*, 3(4):3355–3362, 2018. 1, 2, 4
- [6] Xinke Deng, Yu Xiang, Arsalan Mousavian, Clemens Eppner, Timothy Bretl, and Dieter Fox. Self-supervised 6d object pose estimation for robot manipulation. *arXiv preprint arXiv:1909.10159*, 2019. 2
- [7] Amaury Depierre, Emmanuel Dellandréa, and Liming Chen. Jacquard: A large scale dataset for robotic grasp detection. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3511–3516. IEEE, 2018. 1, 2, 4, 5
- [8] Kuan Fang, Yuke Zhu, Animesh Garg, Andrey Kurenkov, Viraj Mehta, Li Fei-Fei, and Silvio Savarese. Learning task-oriented grasping for tool manipulation from simulated self-supervision. *arXiv preprint arXiv:1806.09266*, 2018. 1
- [9] Di Guo, Fuchun Sun, Huaping Liu, Tao Kong, Bin Fang, and Ning Xi. A hybrid deep architecture for robotic grasp detection. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1609–1614. IEEE, 2017. 2
- [10] Stefan Hinterstoisser, Vincent Lepetit, Slobodan Ilic, Stefan Holzer, Gary Bradski, Kurt Konolige, and Nassir Navab. Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes. In *Asian conference on computer vision*, pages 548–562. Springer, 2012. 1
- [11] Yun Jiang, Stephen Moseson, and Ashutosh Saxena. Efficient grasping from rgb-d images: Learning using a new rect-angle representation. In *2011 IEEE International Conference on Robotics and Automation*, pages 3304–3311. IEEE, 2011. 1, 2, 4, 5
- [12] Wadim Kehl, Fabian Manhardt, Federico Tombari, Slobodan Ilic, and Nassir Navab. Ssd-6d: Making rgb-based 3d detection and 6d pose estimation great again. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1521–1529, 2017. 2
- [13] Ian Lenz, Honglak Lee, and Ashutosh Saxena. Deep learning for detecting robotic grasps. *The International Journal of Robotics Research*, 34(4-5):705–724, 2015. 2
- [14] Sergey Levine, Peter Pastor, Alex Krizhevsky, Julian Ibarz, and Deirdre Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International Journal of Robotics Research*, 37(4-5):421–436, 2018. 2, 4
- [15] Hongzhuo Liang, Xiaojuan Ma, Shuang Li, Michael Görner, Song Tang, Bin Fang, Fuchun Sun, and Jianwei Zhang. Pointnetgpd: Detecting grasp configurations from point sets. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 3629–3635. IEEE, 2019. 2, 4
- [16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision (ECCV)*, pages 740–755. Springer, 2014. 4
- [17] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio Ojea, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. *arXiv preprint arXiv:1703.09312*, 2017. 1, 2, 4
- [18] Douglas Morrison, Peter Corke, and Jürgen Leitner. Closing the loop for robotic grasping: A real-time, generative grasp synthesis approach. *arXiv preprint arXiv:1804.05172*, 2018. 2
- [19] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof graspnet: Variational grasp generation for object manipulation. *arXiv preprint arXiv:1905.10520*, 2019. 2
- [20] Van-Duc Nguyen. Constructing force-closure grasps. *The International Journal of Robotics Research*, 7(3):3–16, 1988. 2, 4, 5
- [21] Lerrel Pinto and Abhinav Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In *2016 IEEE international conference on robotics and automation (ICRA)*, pages 3406–3413. IEEE, 2016. 1, 2, 4
- [22] Yuzhe Qin, Rui Chen, Hao Zhu, Meng Song, Jing Xu, and Hao Su. S4g: Amodal single-view single-shot se (3) grasp detection in cluttered scenes. *arXiv preprint arXiv:1910.14218*, 2019. 2
- [23] Joseph Redmon and Anelia Angelova. Real-time grasp detection using convolutional neural networks. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1316–1322. IEEE, 2015. 1, 2
- [24] Andreas ten Pas, Marcus Gualtieri, Kate Saenko, and Robert Platt. Grasp pose detection in point clouds. *The International Journal of Robotics Research*, 36(13-14):1455–1473, 2017. 1, 2, 4
- [25] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. *arXiv preprint arXiv:1711.00199*, 2017. 2, 4
- [26] Xinchen Yan, Jasmined Hsu, Mohammad Khansari, Yunfei Bai, Arkanath Pathak, Abhinav Gupta, James Davidson, and Honglak Lee. Learning 6-dof grasping interaction via deep geometry-aware 3d representations. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1–9. IEEE, 2018. 1, 2, 4
- [27] Hanbo Zhang, Xuguang Lan, Site Bai, Xinwen Zhou, Zhiqiang Tian, and Nanning Zheng. Roi-based robotic grasp

detection for object overlapping scenes. *arXiv preprint arXiv:1808.10313*, 2018. [1](#), [2](#), [4](#)

- [28] Zelin Zhao, Gao Peng, Haoyu Wang, Hao-Shu Fang, Chengkun Li, and Cewu Lu. Estimating 6d pose from localizing designated surface keypoints. *arXiv preprint arXiv:1812.01387*, 2018. [2](#)