Efficient End-to-End 6-Dof Grasp Detection Framework for Edge Devices with Hierarchical Heatmaps and Feature Propagation

Kaiqin Yang[†], Yixiang Dai[†], Guijin Wang, Siang Chen[™]

**Department of Electronic Engineering

**Tsinghua University*

csa21@mails.tsinghua.edu.cn

†Equal Contribution.

Corresponding author.

Abstract-6-DoF grasp detection is critically important for the advancement of intelligent embodied systems, as it provides feasible robot poses for object grasping. Various methods have been proposed to detect 6-DoF grasps through the extraction of 3D geometric features from RGBD or point cloud data. However, most of these approaches encounter challenges during real robot deployment due to their significant computational demands, which can be particularly problematic for mobile robot platforms, especially those reliant on edge computing devices. This paper presents an Efficient End-to-End Grasp Detection Network (E3GNet) for 6-DoF grasp detection utilizing hierarchical heatmap representations. E3GNet effectively identifies highquality and diverse grasps in cluttered real-world environments. Benefiting from our end-to-end methodology and efficient network design, our approach surpasses previous methods in model inference efficiency and achieves real-time 6-Dof grasp detection on edge devices. Furthermore, real-world experiments validate the effectiveness of our method, achieving a satisfactory 94% object grasping success rate.

Index Terms—Robotics, 6-Dof Grasp Detection, Feature Pyramid Network, Heatmap, Feature Propagation, Edge Devices

I. INTRODUCTION

Visual-guided robotic grasping is fundamental in embodied intelligence. Despite the significant progress made in single object grasping, it remains an challenge to grasp cluttered objects efficiently and precisely in unstructured environments.

Traditional grasping methods utilized model-based prior with full knowledge of objects' physical models [1], which is challenging to use in real world. Advances in deep learning have driven the development of model-free grasp detection methods based on visual information in open-world scenes. Earlier methods generate 4-DoF planar grasps from a single RGBD observation [2], [3]. Despite its efficiency, these representations constrain the gripper perpendicular to the camera plane, which sacrifices the degree of freedom and limits the performance in complex scenes.

Recent advancements in 6-DoF grasp detection have garnered significant attention due to their wide-ranging applications, allowing robots to grasp objects from various orientations. Previous research has proposed methods based on sampling and evaluation to generate 6-DoF grasp poses by sampling and assessing grasp candidates from point clouds

[4], [5]. Other approaches have employed deep neural networks to directly regress relevant attributes in an end-to-end manner [6]. However, both methods could be more efficient and accurate. More recent studies focus on identifying high graspable regions by incorporating global information from scene point clouds or RGBD images [7]–[9]. Though effective, these methods primarily utilize 3D feature encoders, which entail higher computational costs and may not be suitable for edge-devices deployment.

In this work, we propose a novel Efficient End-to-End 6-Dof Grasp detection Network (E3GNet) based on hierarchical heatmap representations, which succeed in detecting highquality and diverse grasps for real-world object clutters. To the best of our knowledge, our method is the first to achieve realtime 6-Dof grasp detection on edge devices. Our framework first constructs a Global Location Heatmap FPN (feature pyramid network) with a lightweight encoder, Geometry-aware MobileOne, to obtain multi-scale features and predict grasp location heatmaps. Then, the Region Feature Propagation module aggregates graspable region features under the guidance of the location heatmaps. Finally, the graspable region features are fed into a specially designed Rotation Heatmap generation model for grasp rotation detection and refinement, forming scene-level 6-Dof grasps. Experiments on the large-scale grasp dataset show that E3GNet detects more precise and diverse grasps than previous methods on the benchmark. Model inference efficiency experiments on multiple platforms, including edge devices, further prove the computation efficiency of our method. Real-world grasp experiments on the real robot also yield a satisfactory 94% grasp success rate.

The detailed contributions of this work are as follows.

- We propose a novel efficient end-to-end 6-Dof grasp detection framework (E3GNet), realizing real-time 6-Dof grasp detection on edge devices.
- We design a novel Region Feature Propagation module and a Rotation-Heatmap-Based Grasp Detection technique to achieve efficient and precise grasp detection.
- We develop a Global Location Heatmap FPN combined with a lightweight encoder, Geometry-aware MobileOne to efficiently obtain multi-scale features and locate grasps.

II. RELATED WORK

A. 6-Dof Grasp Detection

With advancements in deep learning and the availability of large-scale grasp detection datasets, the performance of grasp detection models has seen substantial improvement. S4G [10] extracts features from single-view point cloud and regress grasp poses and grasp quality scores. REGNet [8] builds a three-stage network to perform three tasks: sampling grasp points, generating grasp candidates, and optimizing grasp poses. Wang et al. [7] propose GSnet, which samples seed points from the scene point cloud and aggregates features in the cylindrical area around them to generate grasping configurations. Tang et al. [11] apply multi-radius cylindrical sampling to fuse local features and refine the distance between grasping candidate poses by upsampling. Dai et al. [12] propose the GraspNeRF framework to detect 6-DoF grasp poses of transparent objects based on multiple RGB images. Ma et al. [13] use domain prior knowledge to enhance the generalization ability of the grasp detection network. In contrast to existing literature, our research delves into more effective methods for feature extraction, introducing hierarchical heatmaps to improve both the quality and efficiency of grasp detection.

B. Efficient Grasp Detection

The emergence of large-scale grasping datasets has established a robust data foundation for developing high-precision models. However, this advancement also brings increased training time, model size, and resource consumption. Nie et al. [14] present a lightweight 4-Dof grasping detection network that enhances performance and generalization through knowledge distillation. However, suction grasping and 4-DOF grasping are influenced by the structures of objects and scenes, making it challenging to achieve high-quality grasping, especially for irregular objects. HGGD [9] utilizes heatmap to focus on the graspable region and speed up grasp detection, but it still relies on the 3D PointNet [15] and feature fusion to further detect grasps. Wu et al. [16] have designed an economic grasp detection framework that effectively reduces training resource consumption while maintaining grasping performance. Building on the concept of generative synthesis, Wong et al. [17] introduce the Fast GraspNeXt framework, which efficiently handles multiple tasks, including the detection of occluded objects and the generation of suction grasp heatmaps. Different from the above works, we design a lightweight end-to-end framework with region feature propagation, achieving realtime inference capabilities on edge devices.

III. PROBLEM STATEMENT

Our task is to efficiently generate a set of grasp poses G from a given single-view RGBD image $I \in R^{H \times W \times 4}$ and camera intrinsics. A 6-Dof grasp configuration $g \in G$ can be defined as:

$$\mathbf{g} = (x, y, z, \theta, \gamma, \beta, w),$$

where (x,y,z) is the 3-Dof location of the grasp pose and $(\theta,\gamma,\beta)\in[-\frac{\pi}{2},\frac{\pi}{2}]$ represents the 3-Dof rotation of the grasp

pose in the form of Euler angles in the gripper coordinate system. In addition, an extra parameter w is predicted, which refers to the parallel-jaw width.

IV. METHOD

Our efficient 6-DoF grasp detection framework consists of three primary components: the Global Location Heatmap FPN, Region Feature Propagation, and Regional Rotation-Heatmap-based Grasp Detection. As demonstrated in Fig. 1, the core innovation of our framework is its capability to effectively extract multi-scale features using a lightweight encoder and the Heatmap FPN, which are subsequently propagated to the Regional Rotation-Heatmap-based Grasp Detection stage, enabling precise predictions of grasp rotations. The overall methodology significantly reduces computational complexity while enhancing the quality and robustness of grasp detection.

A. Global Location Heatmap FPN

For feature extraction in the dense image/3D space, differing from the heavy 3D encoder in previous methods [7], [18], [19], we design our encoder in a more efficient fully-convolutional way. We draw inspiration from the RGB encoder MobileOne [20] designed for mobile phones to create a lightweight encoder for RGBD feature extraction. However, the standard MobileOne is inadequate for processing RGBD images due to the loss of geometric information. To preserve crucial 3D geometry, we integrate camera intrinsics into the network by generating an additional positional mesh grid for each pixel in the RGBD image with the camera imaging model:

$$x = \frac{u - c_x}{1000 \times f_x}, \quad y = \frac{v - c_y}{1000 \times f_y}$$

where (u,v) are the pixel coordinates and (f_x, f_y, c_x, c_y) are the camera intrinsics. The encoder, named Geometry-aware MobileOne, takes the RGBD images and the mesh grid as a 6-channel input to extract semantic and geometry features.

Moreover, recognizing the variability in grasp widths in real-world scenarios, we incorporate the widely used feature pyramid network (FPN) [21] to develop multi-scale scene features from the encoder's multi-layer outputs. Building upon the effective location heatmap concept from [9], we design an additional heatmap prediction head after the FPN to leverage features and facilitate subsequent feature propagation. In summary, the Global Location Heatmap FPN, featuring the Geometry-aware MobileOne, effectively extracts multi-scale features, identifies graspable regions, and is optimized for deployment on edge devices.

B. Region Feature Propagation

Though more efficient than previous work, the two-stage framework introduced in [9] exhibits a linearly increasing computational cost associated with the cropping of regions. Differing from HGGD, our approach propagates region-wise features from the FPN under the guidance of the grasp location heatmaps, thus reducing extra computation overhead via feature sharing. Instead of directly conducting feature pooling [22] on the feature maps, we introduce the grid sampling

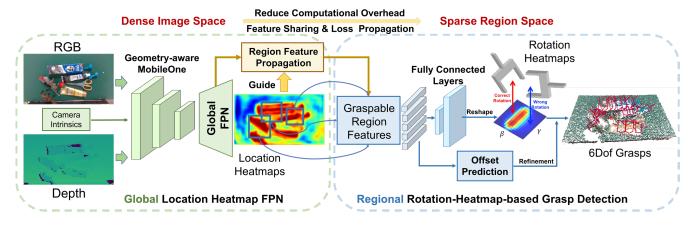


Fig. 1. Overview of the proposed Efficient End-to-End 6-Dof Grasp Detection Framework.

method proposed in [23] to obtain regional feature maps of multiple graspable regions with the same size s. First, we use Farthest Point Sampling [24] to identify the region centers c in the location heatmaps with high graspability values. Then, the local region grid indexes are generated according to the normalized depth values d of the centers:

$$f_{region} = \operatorname{gridsample}(f_{scene}, c + d \times \operatorname{meshgrid}(-\frac{s}{2}, \frac{s}{2})),$$

which means the region feature maps are centered at c with size $s \times s$. To preserve the features' continuity, we use bilinear interpolation in grid sampling. In contrast to a fixed cropping radius [8], [9], we introduce dynamic-scale region cropping by random sampling region size s from a uniform distribution during model training to force our method to learn more robust feature representation. This process can also be applied during inference as a form of test-time augmentation to enhance performance. However, for efficiency, we typically conduct the grid sampling and the feature propagation process just once.

C. Regional Rotation-Heatmap-based Grasp Detection

As proved in pioneer work [7]–[9], [18], anchor-based grasp detectors are efficient for grasp rotation prediction. Following [9], we construct a series of uniform rotation anchors for (γ, β) , whose Cartesian product forms a 2D rotation heatmap. Then, grasp detection can be framed as a semantic segmentation problem within the rotation heatmap. Instead of directly applying anchor-based classification, we enhance the grasp rotation precision by developing higher-resolution rotation heatmaps. In our implementation, to exploit the relationship between different rotations, we employ two fully connected layers to output the rotation heatmaps. Additionally, we incorporate a refinement module to produce grasp widths and additional center offsets, improving grasp detection quality.

V. EXPERIMENTS

A. Dataset, Metrics and Results

We utilize the GraspNet-1Billion dataset [18] to train and evaluate our network, which comprises 100 real-world scenes, each containing RGBD images captured from 256 viewpoints. The test split contains 90 scenes, which are categorized into

three sub-datasets: seen, similar, and novel, depending on the type of objects. We adopt the metric AP proposed by the GraspNet-1Billion dataset, which is the average of force-closure accuracy [26], AP $_{\mu}$, when the friction coefficient μ ranges from 0.2 to 1, with an interval of 0.2. For each scene, AP $_{\mu}$ represents the average grasp precision of the top 50 grasps after non-maximum suppression to evaluate grasp diversity.

We evaluate the overall performance of E3GNet against several grasp detection algorithms, all trained on the same dataset split. As indicated in Table I, E3GNet demonstrates impressive results, achieving an average of 52.18 mAP across all the test scenes. Compared to other approaches, E3GNet surpasses the current state-of-the-art on all metrics, particularly excelling in the unseen categories, mainly benefiting from the well-designed feature extraction and propagation scheme.

B. Model Efficiency Experiments

In addition to the quality of grasp detection, the inference speed of the model significantly impacts real-world grasping execution, particularly on certain edge devices. To thoroughly evaluate model efficiency, we tested various models across two types of devices, totaling four different platforms. We selected the one-time inference time—from data input to grasp output—as our evaluation metric. As shown in Table II, E3GNet demonstrates superior performance in network inference speed on all tested platforms, highlighting the effectiveness of our framework design and network optimization.

C. Ablation Studies

Ablation studies are conducted on the three primary components of E3GNet, with HGGD [9] as the baseline. The results presented in Table III underscore the impact of these components. As anticipated, our novel framework, incorporating Region Feature Propagation, significantly improves model inference efficiency by minimizing redundant computations and largely enhancing overall performance. Transitioning to a more efficient lightweight Global Location Heatmap FPN yields faster inference speeds, albeit with a slight decrease in grasp detection quality. Additionally, Rotation-Heatmap-based Grasp Detection moderately improves grasp detection quality without incurring any extra computational costs.

Method	Seen			Similar			Novel			Average
	AP	$AP_{0.8}$	$AP_{0.4}$	AP	$AP_{0.8}$	$AP_{0.4}$	AP	$AP_{0.8}$	$AP_{0.4}$	mAP
GPD [5]	22.87	28.53	12.84	21.33	27.83	9.64	8.24	8.89	2.67	17.48
PointnetGPD [4]	25.96	33.01	15.37	22.68	29.15	10.76	9.23	9.89	2.74	19.29
GraspNet [18]	27.56	33.43	16.95	26.11	34.18	14.23	10.55	11.25	3.98	21.41
TransGrasp [19]	39.81	47.54	36.42	29.32	34.80	25.19	13.83	17.11	7.67	27.65
HGGD [†] [9]	58.35	66.54	55.96	47.93	56.91	41.86	22.10	27.37	14.31	42.79
SBGrasp [25]	58.95	68.18	54.88	52.97	63.24	46.99	22.63	28.53	12.00	44.85
GSNet [7]	65.70	76.25	61.08	53.75	65.04	45.97	23.98	29.93	14.05	47.81
E3GNet	69.56	78.35	66.30	60.84	72.81	52.07	26.74	33.12	15.43	52.38

^{*}The **best scores** are displayed in **bold**.

TABLE II
MODEL INFERENCE EFFICIENCY ON DIFFERENT DEVICE PLATFORMS

Platform	Method Inference Time / ms						
riauoriii	GSNet	GraspNet	HGGD	E3Gnet			
Desktop Devices							
Nvidia RTX 3090 ¹	104.1	85.2	24.3	5.9			
Nvidia RTX 3060Ti ²	154.2	102.0	30.7	9.2			
Edge Devices							
Nvidia Jetson TX2	Failed	424.8	649.4	157.9			
Nvidia Jetson Xavier NX	Failed	363.1	418.1	126.7			

^{*}All time are averaged by 200 trials and shown in miliseconds.

TABLE III
METHOD COMPONENT ABLATION, SHWOING RESULT ON REALSENSE
SPLIT AND TIME TESTED ON THE NVIDIA JETSON TX2 PLATFORM

Components			Seen	Similar	Novel	Time/ms	
RFP	GLH	RRH	Seen	Sillilai	Novei	1 IIIIe/IIIS	
			58.35	47.93	22.10	649.4	
\checkmark			68.20	58.96	26.58	186.4	
\checkmark	\checkmark		67.62	58.72	26.53	158.2	
\checkmark	\checkmark	\checkmark	69.56	60.84	26.74	157.9	

RFP: Region Feature Propagation GLH: Global Location Heatmap FPN

RRH: Regional Rotation-Heatmap-based Grasp Detection

D. Realworld Grasping Result

To validate the effectiveness of E3GNet in real-world settings, we perform real-world robotic grasping experiments on a UR-5e robotic arm equipped with a Robotiq 2F-85 parallel-jaw gripper. A Realsense-D435i camera is mounted to capture egocentric RGBD images. Our real-world experiment employs a collection of 22 objects with various shapes, generating 6 cluttered scenes composed of 7 to 9 randomly selected items in different poses. We focus on object grasping to clear these clutters. The robot conducts grasp detection and executes grasps until no grasps are identified or a maximum of 12 attempts are reached. Following [9], [18], we apply Success Rate (successful grasps / total attempts) and Completion Rate (cleared scenes / total scenes) as evaluation metrics.

Table IV presents the performance of our method in real-

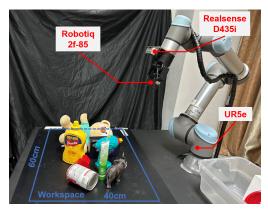


Fig. 2. Real-world robot experiment settings.

TABLE IV
REAL-WORLD CLUTTER CLEARANCE RESULTS

Scene	Objects	E3GNet	HGGD		
1	7	7 / 8	7/7		
2	8	8 / 8	8 / 11		
3	8	8 / 9	7 / 12		
4	8	8 / 9	8 / 11		
5	8	8 / 8	8 / 8		
6	9	9/9	9 / 10		
Success Rate		94% (48 / 51)	80%(47 / 59)		
Comple	etion Rate	100% (6 / 6)	83%(5 / 6)		

world environments. Our approach achieves an impressive average grasp success rate of 94 % and successfully clears all 6 scenes, demonstrating its ability to generate high-quality grasps. However, some failures were noted, particularly when the gripper collided with other objects in the clutters. A demo video of real-world grasping can be found at youtube.

VI. CONLCUSION

In this paper, we introduce an efficient end-to-end framework for 6-Dof grasp detection. By employing a hierarchical heatmap-based design, our method enables real-time detection of high-quality 6-Dof grasps and is suitable for deployment on edge devices. Our framework demonstrates state-of-theart performance on the GraspNet-1Billion Dataset [18]. Additionally, quantitative experiments involving real-world robot grasping validate the effectiveness of our approach.

^{†:} Reimplemented with official codebase.

¹Intel 13900KF CPU, Nvidia RTX 3090 GPU and ubuntu 20.04.

²AMD 5600X CPU, Nvidia RTX 3060Ti GPU and ubuntu 20.04.

REFERENCES

- J. Bohg, A. Morales, T. Asfour, and D. Kragic, "Data-driven grasp synthesis—a survey," *IEEE Transactions on robotics*, vol. 30, no. 2, pp. 289–309, 2013.
- [2] S. Kumra, S. Joshi, and F. Sahin, "Antipodal robotic grasping using generative residual convolutional neural network," in *IEEE/RSJ Interna*tional Conference on Intelligent Robots and Systems (IROS), 2020.
- [3] D. Morrison, P. Corke, and J. Leitner, "Closing the loop for robotic grasping: A real-time, generative grasp synthesis approach," arXiv preprint arXiv:1804.05172, 2018.
- [4] H. Liang, X. Ma, S. Li, M. Görner, S. Tang, B. Fang, F. Sun, and J. Zhang, "Pointnetgpd: Detecting grasp configurations from point sets," in 2019 International Conference on Robotics and Automation (ICRA). IEEE, 2019, pp. 3629–3635.
- [5] A. Ten Pas, M. Gualtieri, K. Saenko, and R. Platt, "Grasp pose detection in point clouds," *The International Journal of Robotics Research*, vol. 36, no. 13-14, pp. 1455–1473, 2017.
- [6] P. Ni, W. Zhang, X. Zhu, and Q. Cao, "Pointnet++ grasping: Learning an end-to-end spatial grasp generation algorithm from sparse point clouds," in 2020 IEEE International Conference on Robotics and Automation (ICRA), 2020, pp. 3619–3625.
- [7] C. Wang, H.-S. Fang, M. Gou, H. Fang, J. Gao, and C. Lu, "Graspness discovery in clutters for fast and accurate grasp detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision* (ICCV), October 2021, pp. 15964–15973.
- [8] B. Zhao, H. Zhang, X. Lan, H. Wang, Z. Tian, and N. Zheng, "Regnet: Region-based grasp network for end-to-end grasp detection in point clouds," in 2021 IEEE International Conference on Robotics and Automation (ICRA), 2021, pp. 13474–13480.
- [9] S. Chen, W. Tang, P. Xie, W. Yang, and G. Wang, "Efficient heatmapguided 6-dof grasp detection in cluttered scenes," *IEEE Robotics and Automation Letters*, vol. 8, no. 8, pp. 4895–4902, 2023.
- [10] Y. Qin, R. Chen, H. Zhu, M. Song, J. Xu, and H. Su, "S4g: Amodal single-view single-shot se (3) grasp detection in cluttered scenes," in *Conference on robot learning*. PMLR, 2020, pp. 53–65.
- [11] W. Tang, K. Tang, B. Zi, S. Qian, and D. Zhang, "High precision 6-dof grasp detection in cluttered scenes based on network optimization and pose propagation," *IEEE Robotics and Automation Letters*, vol. 9, no. 5, pp. 4407–4414, 2024.
- [12] Q. Dai, Y. Zhu, Y. Geng, C. Ruan, J. Zhang, and H. Wang, "Graspnerf: Multiview-based 6-dof grasp detection for transparent and specular objects using generalizable nerf," in 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2023, pp. 1757–1763.
- [13] H. Ma, M. Shi, B. Gao, and D. Huang, "Generalizing 6-dof grasp detection via domain prior knowledge," in *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition (CVPR), June 2024, pp. 18102–18111.
- [14] H. Nie, Z. Zhao, L. Chen, Z. Lu, Z. Li, and J. Yang, "Smaller and faster robotic grasp detection model via knowledge distillation and unequal feature encoding," *IEEE Robotics and Automation Letters*, vol. 9, no. 8, pp. 7206–7213, 2024.
- [15] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [16] X.-M. Wu, J.-F. Cai, J.-J. Jiang, D. Zheng, Y.-L. Wei, and W.-S. Zheng, "An economic framework for 6-dof grasp detection," arXiv preprint arXiv:2407.08366, 2024.
- [17] A. Wong, Y. Wu, S. Abbasi, S. Nair, Y. Chen, and M. J. Shafiee, "Fast graspnext: A fast self-attention neural network architecture for multitask learning in computer vision tasks for robotic grasping on the edge," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2023, pp. 2293–2297.
- [18] H.-S. Fang, C. Wang, M. Gou, and C. Lu, "Graspnet-Ibillion: A large-scale benchmark for general object grasping," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 444–11 453.
- [19] Z. Liu, Z. Chen, S. Xie, and W.-S. Zheng, "Transgrasp: A multi-scale hierarchical point transformer for 7-dof grasp detection," in 2022 International Conference on Robotics and Automation (ICRA). IEEE, 2022, pp. 1533–1539.

- [20] P. K. A. Vasu, J. Gabriel, J. Zhu, O. Tuzel, and A. Ranjan, "Mobileone: An improved one millisecond mobile backbone," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 7907–7917.
- [21] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
- [22] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE transactions on* pattern analysis and machine intelligence, vol. 39, no. 6, pp. 1137–1149, 2016.
- [23] M. Jaderberg, K. Simonyan, A. Zisserman et al., "Spatial transformer networks," Advances in neural information processing systems, vol. 28, 2015
- [24] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," Advances in neural information processing systems, vol. 30, 2017.
- [25] H. Ma and D. Huang, "Towards scale balanced 6-dof grasp detection in cluttered scenes," in *Conference on robot learning*. PMLR, 2023, pp. 2004–2013.
- [26] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan, X. Liu, J. A. Ojea, and K. 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.