

Hybrid Physical Metric For 6-DoF Grasp Pose Detection

Yuhao Lu, Beixing Deng, Zhenyu Wang, Peiyuan Zhi, Yali Li, Shengjin Wang
Tsinghua University

Abstract—6-DoF grasp pose detection of multi-grasp and multi-object is a challenge task in the field of intelligent robot. To imitate human reasoning ability for grasping objects, data driven methods are widely studied. With the introduction of large-scale datasets, we discover that a single physical metric usually generates several discrete levels of grasp confidence scores, which cannot finely distinguish millions of grasp poses and leads to inaccurate prediction results. In this paper, we propose a hybrid physical metric to solve this evaluation insufficiency. First, we define a novel metric is based on the force-closure metric, supplemented by the measurement of the object flatness, gravity and collision. Second, we leverage this hybrid physical metric to generate elaborate confidence scores. Third, to learn the new confidence scores effectively, we design a multi-resolution network called Flatness Gravity Collision GraspNet (FGC-GraspNet). FGC-GraspNet proposes a multi-resolution features learning architecture for multiple tasks and introduces a new joint loss function that enhances the average precision of the grasp detection. The network evaluation and adequate real robot experiments demonstrate the effectiveness of our hybrid physical metric and FGC-GraspNet. Our method achieves 90.5% success rate in real-world cluttered scenes. Our code is available at <https://github.com/luyh20/FGC-GraspNet>.

I. INTRODUCTION

Grasping is one of the most fundamental and important tasks in the field of robotic manipulation. Recently data-driven methods [11], [24], [29] have been developed to reason robust grasps under various settings. For the point cloud data of observed scenes, the model outputs grasp poses by a deep neural network [3], [28], [36]. However, generating reliable and humanoid grasps for multiple objects in unstructured and cluttered environments is still a challenge.

The performance and efficacy of data-driven methods mainly depends on two aspects: inferring the grasp quality and predicting the grasp pose. Since it is laborious and costly to annotate the potential success rates of large-scale 6-DoF grasp poses by human, a solid grasp evaluation mechanism to infer the grasp quality is of vital importance. Some recent grasp pose detection methods [3], [11], [14], [29], [36] apply the physics analytic approaches to evaluate the quality of grasp poses. Thereinto, the force-closure metric [30] is a mainstream evaluation metric. Although it has been utilized by many works [3], [14], it is still limited. The main problem is that it only provides a binary outcome under a coefficient of friction, and usually generates confidence scores in several ranked discrete levels, like ten bins in [3]. Such discrete score levels necessarily exist noise. As is illustrated in the top of Fig. 1, some grasp poses apparently achieve different grasp performances but are assigned in the same level score under force-closure metric. Because of the noise, the network is likely to predict confidence scores deviated from the

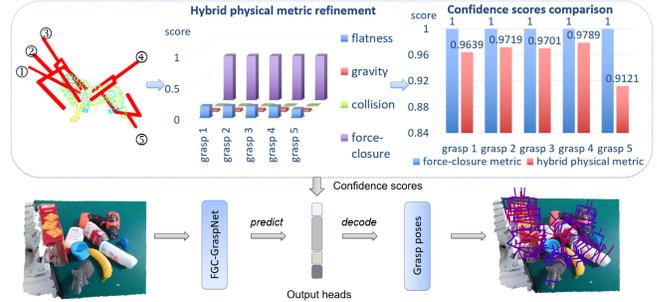


Fig. 1. Top: Grasp confidence scores refinement. There are five grasp candidates labeled in the lion model. The middle chart shows the composition of hybrid physical metric. The right chart is the comparison results. Bottom: Grasp pose detection pipeline. The input data forwards through our FGC-GraspNet and predicts grasp poses.

truth. Such a confidence score gap between real world grasp success rates and calculated confidence scores affects the adaptability of the network in real world application. In addition, the force-closure metric only describes a single physical characteristic, which is quite insufficient to reason reliable grasp poses when encountering novel objects and degrades the robustness of grasp pose detection.

To tackle these issues above, in this paper, we propose a hybrid physical metric to solve the evaluation insufficiency. To reduce the influence of label noise, we leverage more comprehensive physical information to refine grasp confidence scores. According to human grasp habits, we pay attention to the contact points between the two-finger gripper and the corresponding object. The first evaluation metric we adopt is the flatness metric, which aims to measure the flatness of contact points. The second is the gravity center metric, and the motivation is to balance the gravity and the grasping pressure. The third is the collision perturbation metric, which prevents possible collisions between grasp end points and contact points. Together with force-closure metric, we present the hybrid physical metric to evaluate the final grasp confidence scores. As is shown in the top part of Fig. 1, our designed hybrid physical metric is used to distinguish grasp candidates more reasonably than force-closure metric.

Meanwhile, to learn this fine hybrid metric more effectively, we further design a multi-resolution network called FGC-GraspNet. During the training process, the grasp confidence scores are passed into different loss functions for the multi-task learning paradigm [25], [26], such as the approaching direction or the in-plane rotation prediction [3]. However, different tasks usually require different kinds of information. Specifically, the predecessor tasks, including

foreground segmentation and approaching vector prediction, often need scene-level and object-level information, while the post tasks like in-plane rotation and depth prediction focus more on object-level and point-level information. Hence, FGC-GraspNet is a multi-resolution network. Based on the hierarchical feature learning of PointNet++ [22], the features of low-resolution point sets are extracted for the predecessor tasks. The features of high-resolution point sets forward a local attention module to gather region information for the post tasks. Furthermore, to better accommodate to our hybrid metric, we adopt a new joint learning loss function. This loss regresses the confidence scores more sufficiently and contains a grasp depth classification loss.

We implement the hybrid physical metric on the GraspNet-1Billion dataset [3], and obtain more elaborate and physically meaningful grasp confidence scores. We then use the new grasp confidence scores in our grasp pose detection work. Experiment results demonstrate that our FGC-GraspNet under new grasp confidence scores achieves significantly better performance.

To summarize, our main contributions are as follows:

- We propose a hybrid physical metric to refine the grasp confidence scores. This metric adopts more comprehensive physical descriptions thus is more reasonable.
- We design a multi-resolution network FGC-GraspNet. By utilizing information of different levels of resolution for multiple tasks and adopting new joint loss function, our network better adapts to our hybrid physical metric.
- Extensive experiments show that the hybrid physical metric is beneficial for increasing the success rates in reality, and our network significantly promotes the ability of grasp prediction.

II. RELATED WORK

A. 6-DoF Grasp Pose Detection

Most of recent 6-DoF grasp pose detection works are based on data-driven methods, where 6-DoF is decoupled to 3D position and 3D rotation vector for the movement of the robotic arm. Compared to rectangle grasp representation [4], [8], 6-DoF allows multi-view camera inputs and provides more possible approaching directions of grasp poses. Many recent methods [5], [15], [16], [23], [29], [31], [36] attempt to learn predicting grasp poses based on deep learning. PointNetGPD [14] samples grasp candidates and evaluates the grasp quality based on the network PointNet [21]. GraspNet-1Billion [3] builds a large-scale grasp dataset and proposes a baseline method for learning grasp poses. REGNet [36] use group region features to predict grasp proposals. Contact-GraspNet [28] proposes a new grasp pose representation and implement it in [2]. A object instance segmentation network [34] is utilized in [17] to tackle the problem of grasping objects in a cluttered scene. These methods greatly enrich the solutions of grasp pose detection. However, most of them evaluate the grasp quality at discrete levels. Therefore, we propose hybrid physical metric and design a multi-resolution network to learn the refined confidence scores.

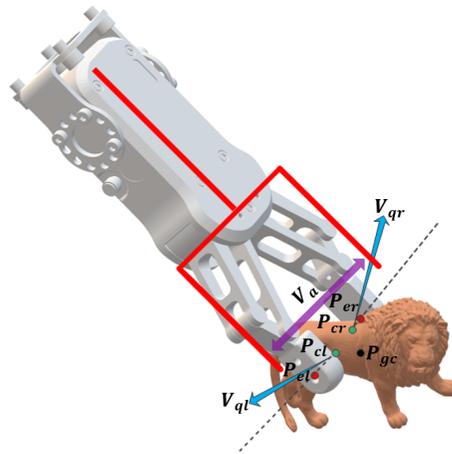


Fig. 2. Grasp pose demonstration. P_{cl} , P_{cr} are the contact points in the object, P_{el} , P_{er} are the end points of gripper, P_{gc} is the gravity center of the object, V_a is the antipodal vector, V_{ql} , V_{qr} are the normals of the contact points. The dotted line is the reference line of V_n . The gripper model is from RG2¹. The lion model is from the dataset [3].

B. Grasp Evaluation Metrics

Grasp quality evaluation is essential component of supervised learning grasp detection methods. A reliable grasp labels should reflect the grasping success rate in real robotic manipulation accurately. A recent work like [12] produces grasp labels through a large-scale real robotic grasping attempts. Some works [1], [2], [9], [15], [35] utilize physics simulators to annotate grasp confidence scores. Many analysis approaches [3], [11], [14], [29], [36] generate confidence scores by calculating the contact physics about the geometry of gripper configuration and object mesh models. Defining a reasonable grasp quality metric is an open problem. Fingertip space is proposed in [7] to search stable grasps by considering both the local geometry of object surface and the fingertip geometry and a flatness criteria is defined to filter points. PointNetGPD [14] modify the coefficient of friction to get a discrete score based on force-closure [18] and use grasp wrench space (GWS) [10] analysis as a auxiliary. GraspNet-1Billion [3] also calculate friction coefficient of grasp poses at 10 bins. REGNet [36] measure the angle between the force direction and contacts' normals. Four different grasp stability metrics are applied in [9], including ϵ -metric [20] and so on. These metrics have respective limitations for 6-DoF grasp pose detection. Hence, we propose the hybrid physical metric in this paper.

C. Point Cloud Learning Methods

Point cloud learning is one of necessary module during grasp pose detection methods of data driven. A variety of approaches are proposed for the features learning of point cloud, consisting of point-wise MLP methods such as PointNet [21], PointNet++ [22], convolution-based methods such as PointConv [33], PointCNN [13], and some newly proposed transformer-based methods such as PCT [6], PT [37]. In the field of 3D object detection, Local-Global transformer module [19] and a two-stage architecture for

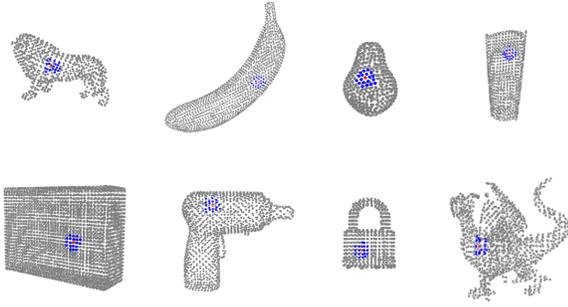


Fig. 3. Flatness measurement of partial object mesh models. These object models are from the dataset [3]. The red point is the flattest point, indicating the highest S_f among all points. The blue points are the neighbors points.

multi-task learning [27] are proved useful in generating 3D object proposals. Currently, PointNet++ [22] is used widely in the area of grasp pose detection.

III. HYBRID PHYSICAL METRIC

In this section, we introduce our proposed hybrid physical metric. Given the 3D object mesh model and the grasp pose annotated in the model, the grasp evaluation process aims to predict the confidence scores for different grasp poses. We use the same 6-DoF grasp representation to define the grasp pose with the dataset [3]. Based on the previous force-closure metric, we further adopt flatness metric, gravity center metric and collision perturbation metric, to generate more accurate grasp confidence scores on the GraspNet 1-Billion dataset [3]. Specifically, we focus on the contact points, the critical geometric points and vectors, which are displayed in Fig. 2.

Flatness Metric. For the grasp action of a two-finger gripper, intuitively, the grasp quality is higher when the contact region is more flat. To utilize this information, we quantify the flatness of the contact region in two steps.

First, we calculate the flatness of points in 3D mesh models. The similarity of the local normal vectors in the query point region can be used to measure the flatness score of the points, denoted with S_{f1} . We use cosine distance between k-nearest neighbors' normals and the query point normal to calculate it. Partial flatness measurement results with our designed score are shown in Fig. 3. Second, we consider the perpendicularity between the antipodal direction and the contact region. We calculate the consistency between the antipodal vector and the contact point normal with the cosine distance as the consistency score S_{f2} . The final score S_f is obtained through the multiplication operation of these two scores. The specific calculation operation can be formulated as follows:

$$\begin{aligned}
 S_{f1} &= \frac{1}{2K} \sum_{q=1}^2 \sum_{n=1}^K \frac{\langle V_q, V_n \rangle}{\|V_q\| \cdot \|V_n\|} \\
 S_{f2} &= \frac{1}{2} \sum_{q=1}^2 \frac{|\langle V_a, V_q \rangle|}{\|V_a\| \cdot \|V_q\|} \\
 S_f &= S_{f1} \cdot S_{f2}
 \end{aligned} \tag{1}$$

where V_q is the normal vector of the query point, including two contact points, V_n is the normal vector of its neighbor points, K is the number of neighbors and V_a is the antipodal vector. $\langle \cdot, \cdot \rangle$ denotes the inner product operation.

Gravity Center Metric. Considering the fact that the grasp candidate whose antipodal force is closer to the gravity center of object is more steady, we propose the gravity center metric. We thus quantify this mechanical relationship between the gravity and the pressure into a distance metric. In the light of the geometry of the grasp pose, the two contact points are linked to an antipodal line. We adopt the euclidean distance between the gravity center point and the antipodal line as the gravity score S_g :

$$S_g = \frac{\|(P_{cl} - P_{gc}) \times (P_{cr} - P_{gc})\|}{\|P_{cl} - P_{cr}\|} \tag{2}$$

where P_{cl} , P_{cr} and P_{gc} are coordinates of the left contact point, right contact point and the gravity center point. The S_g is normalized and converted to $1 - S_g$ in practice.

Collision Perturbation Metric. We observe that the grasp candidates are prone to collision when the end point is close to the object model in real world experiments. Hence, to avoid too close contact, the minimum value of the euclidean distances between the end points and the object contact points is formulated as the collision perturbation score S_c ,

$$S_c = \min(\|P_{el} - P_{cl}\|, \|P_{er} - P_{cr}\|) \tag{3}$$

where P_{el} , P_{er} are coordinates of the left end point and the right end point. The S_c is also normalized.

Hybrid Physical Metric. The hybrid physical metric is a combination of our proposed metrics above and previous force-closure metric. Since different physical viewpoints are adopted, our proposed metric possesses better generalization ability of grasping reasoning. The force-closure metric generates a ten-bin grasp confidence score S_t . The final grasp score S is computed as follows.

$$S = \lambda_t \cdot S_t + \lambda_f \cdot S_f + \lambda_g \cdot S_g + \lambda_c \cdot S_c \tag{4}$$

We set $\lambda_t, \lambda_f, \lambda_g, \lambda_c = 0.7, 0.2, 0.05, 0.05$ in practice.

IV. FGC-GRASPNET

A. Overview

To detect grasp poses in scene-level point clouds, each grasp element constituting the grasp representation needs to be predicted by the network. According to [3], the foreground segmentation needs to be carried out at first, then the grasp representation prediction is decoupled into multiple tasks including the prediction of depth, width, approaching vector and the in-plane rotation. In predecessor tasks, the foreground segmentation considers the overall geometry structure and the grasp approaching vector prediction is closely associated with the direction to the ground in the world space. In comparison, in post tasks, the in-plane rotation, depth and width are usually related to the local

¹https://github.com/ekorudiawan/rg2_simulation

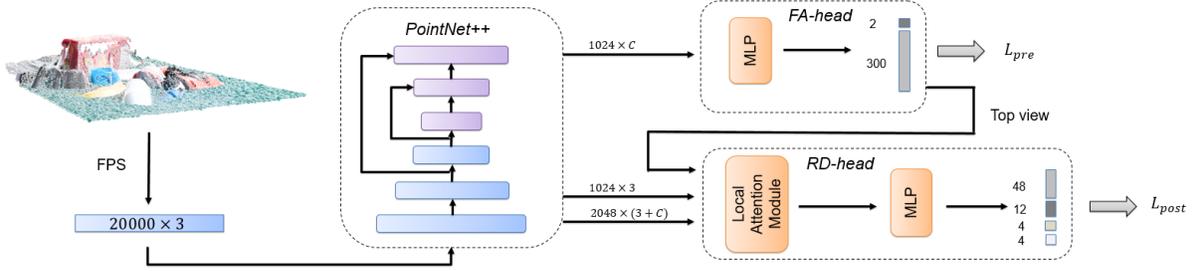


Fig. 4. The architecture of FGC-GraspNet. The input point clouds are sampled by farthest point sampling (FPS) [22] to 20000×3 . The network consists of PointNet++ [22], FA-head (Foreground-Approach-head), RD-head (Rotation-Depth-head), and more details in text.

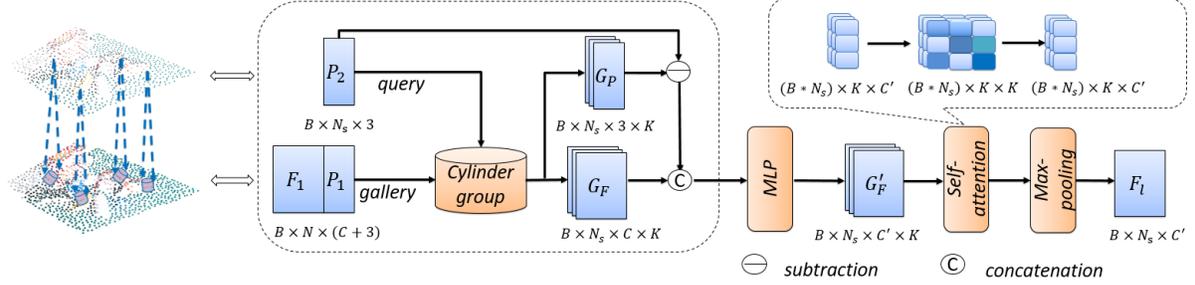


Fig. 5. Local attention module. Left: the diagram shows the process of cylindrical region query between cross-resolution point clouds; Middle: the detailed architecture of cylinder grouping, P_1 , P_2 are respectively high and low resolution point sets from PointNet++ backbone, N_s is the number of seed points, K is the sample number during cylinder grouping; Right: local attention operation outputs final features map F_t , the top part of self-attention unit indicates the dimension of attention map is $K \times K$.

geometry of a single object. Therefore, we design a multi-resolution architecture to extract features for these multiple tasks. As is shown in Fig. 4, our network consists of a base backbone PointNet++ [22] and two branches called FA-head and RD-head. PointNet++ [22] is leveraged to extract features of hierarchical point sets. Features of the low-resolution seed point set pass into the FA-head for the foreground segmentation and point-wise approach direction score regression, while features of high-resolution point set are served for post tasks in the RD-head through our designed local attention module. We finally adopt a new joint learning loss. On one hand, considering our finer confidence score, we perform the regression operation for all the predicted scores instead of only regressing their maximal values. In this way, the network can be supervised more sufficiently. On other hand, since the supervision information for the maximal predicted score is weakened, we add a loss for the depth prediction task. Our new designed loss can be better adapted to our hybrid physical metric.

B. Local Attention Module

The in-plane rotation and depth prediction usually depend on the local geometry structure of the single object model. To extract more abundant and complex local information, we design a local attention module. This module queries region information in high-resolution feature map and evolves it by the self-attention unit [32]. Its structure is shown in the Fig. 5. First, cylinder region query is employed to search the neighbor points of the seed point set P_2 in the high-resolution point set P_1 . More cylinder grouping details can be referred to [3]. This grouping process outputs

point features $G_F \in \mathbb{R}^{B \times N_s \times C \times K}$ and point coordinates $G_P \in \mathbb{R}^{B \times N_s \times 3 \times K}$. Then, G_F is concatenated with the coordinate offsets between query points P_2 and group points G_P . This concatenated feature will be processed by a self-attention layer to enhance the local region attention. The self-attention layer only focuses on capturing region-range contextual information, and its details can be referred to [6], [37]. These group features finally forward through the max pooling layer along the K -dimension to retain the most salient features. Our local attention module integrates the features and coordinates of spatial neighbor points, which better adapts to the in-plane rotation and depth classification tasks.

C. Loss function

For the predecessor tasks in FA-head, we adopt a classification loss to learn the object mask, and a regression loss for view scores to supervise the approach direction learning. The loss L_{pre} is as follow:

$$L_{pre} = \frac{1}{N} \sum_{i=1}^N L_{cls}(\hat{m}_i, m_i) + \frac{1}{N_{reg}} \sum_{i=1}^N \sum_{j=1}^V m_i \cdot L_{reg}(\hat{s}_{ij}, s_{ij}) \quad (5)$$

where m_i is a binary label that it is assigned 1 if the point is of objects, s_{ij} is the view score label used the maximum grasp score in each approach direction, \hat{m}_i , \hat{s}_{ij} represent the corresponding predicted values. $N = 1024$ and $V = 300$.

For the post tasks in RD-head, we regress all 48 grasp scores which are correspond to 48 grasp proposals composed of 12 types of rotation and 4 types of depth. In addition, we predict the rotation direction and depth through a classifica-

TABLE I

EVALUATION FOR DIFFERENT MODELS. THE TABLE SHOWS THE RESULTS ON DATA CAPTURED BY REALSENSE/KINECT RESPECTIVELY

Models	Seen			Unseen			Novel		
	mAP	AP _{0.3}	AP _{0.7}	mAP	AP _{0.3}	AP _{0.7}	mAP	AP _{0.3}	AP _{0.7}
Fang et al. [3]	37.40/32.79	34.75/27.79	18.79/14.75	35.01/30.45	30.25/24.34	17.54/11.18	23.22/21.05	12.36/9.29	3.21/2.48
Ours(no depth)	41.88/37.53	41.59/34.65	24.11/20.11	36.75/31.41	33.47/25.66	19.21/12.18	24.48/21.81	14.61/10.22	3.89/2.61
Ours	49.68/41.09	53.06/40.18	33.73/23.58	40.09/33.05	38.40/28.35	23.31/13.64	26.01/23.27	17.37/12.31	5.03/3.35

TABLE II

RESULTS OF SINGLE OBJECT GRASPING TESTS. FCM IS FORCE-CLOSURE METRIC, HPM IS HYBRID PHYSICAL METRIC.

Object/ID	Banana/5	Apple/12	Nivea.../42	Hosjam/44	Giraffc/55	Weiquan/57	Darlic../58	Lion/67	All
Attempt	30	30	30	30	30	30	30	30	240
Success	25	25	17	28	24	27	27	24	197
FCM	Success Rate	83.33%	83.33%	56.67%	93.33%	80%	90%	80%	82.08%
HPM	Success	27	25	22	29	26	27	25	209
HPM	Success Rate	90%	83.33%	70%	96.67%	86.67%	93.33%	90%	87.08%

tion loss. The grasp width is regressed along the prediction rotation category. The loss L_{post} is formulated as follows:

$$\begin{aligned}
L_{post} = & \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{N_{reg}} \sum_{j=1}^{A \times D} m_i \cdot L_{reg}(\hat{s}_{ij}, s_{ij}) \right) \\
& + \frac{1}{N_{cls}} \sum_{j=1}^A m_i \cdot L_{cls}(\hat{r}_{ij}, r_j) + \frac{1}{N_{cls}} \sum_{j=1}^D m_i \cdot L_{cls}(\hat{d}_{ij}, d_j) \\
& + \alpha \cdot \frac{1}{N_{reg}} \sum_{i=1}^D m_i \cdot L_{reg}(\hat{w}_{ij}, w_{ij})
\end{aligned} \quad (6)$$

where m_i is the same in L_{pre} , s_{ij} is 48 candidate score labels on the top view direction predicted by the FA-head, r_i is the index of maximum grasp score among 12 in-plane rotation directions and d_i is the index of maximum grasp score among 4 depths, w_{ij} is grasp width of the r_i direction grasp pose, \hat{s}_{ij} , \hat{r}_{ij} , \hat{d}_{ij} , \hat{w}_{ij} represent the corresponding prediction values, and $A = 12$, $D = 4$, we set $\alpha = 0.2$.

Finally, the overall joint learning loss L is;

$$L = L_{pre} + \beta \cdot L_{post} \quad (7)$$

We set $\beta = 0.3$ in practice.

V. EXPERIMENTS

A. Network Evaluation

1) *Evaluation Metric*: We conduct experiments on the GraspNet-1Billion dataset [3], and leverage our new grasp confidence scores according to Eq.4. We use the Average Precision (AP) [3] to evaluate the network. Specifically, the grasp pose non-maximum suppression (NMS) and collision detection are used to filter candidates at first. Then, we extract candidates with the top-50 predicted scores, and query the corresponding real confidence scores. We set up different score thresholds, then calculate the mean AP (mAP) under different score thresholds [0, 0.1, 0.3, 0.5, 0.7, 0.9], which are corresponding to friction coefficient thresholds in [3].

TABLE III

RESULTS OF SCENE GRASPING.

Scene	Object ID	Attempt	Success Rate	
			Fang et al. [3]	Ours
Scene1	5,12,39,44,67	10	90%	92%
Scene2	30,42,53,59,61	10	82%	86%
Scene3	35,36,55,57,63	10	92%	92%
Scene4	37,38,40,48,68	10	84%	92%
All	All	40	87%	90.5%

2) *Comparison Studies*: We compare the performance of networks under the new confidence scores generated by the hybrid physical metric. Other parameter settings of the experiments are all the same. Since the numerical distribution of the grasp confidence scores has changed, the AP results are different from [3]. As is illustrated in Table I, our FGC-GraspNet outperforms the network of [3] on both Realsense camera dataset and Kinect camera dataset. On the seen test set, our FGC-GraspNet outperforms previous method by more than 10% AP. On the unseen or novel test set, the improvement is also quite significant. This demonstrates the effectiveness of the FGC-GraspNet.

3) *Ablation Studies*: We also conduct ablation experiments to analysis the contributions of the depth classification loss in our new joint loss function. The results are listed in Table I. Ours(no depth) indicates that we do not adopt the depth classification loss function in the joint learning loss. Instead, we predict the depth by using different height of cylinder like [3]. We notice that our depth loss generally improves the results by more than 7% on the seen test set. This indicates that our network is greatly improved by the depth classification loss function.

B. Hybrid Physical Metric Evaluation

To demonstrate the efficiency of our hybrid physical metric, we conduct single object grasping experiments in the real robotic environment. We select 8 objects in tests. To ensure the robustness of the experiment, the geometric shapes of these objects are representative. We use the FGC-GraspNet to obtain two models. One is trained in new grasp scores generated from the hybrid physical metric, and

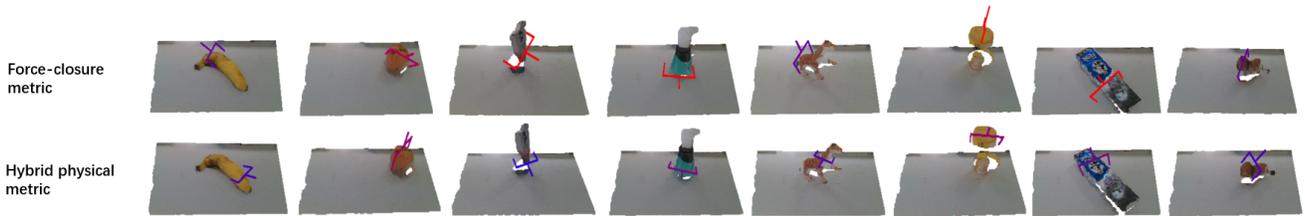


Fig. 6. Single object grasping. Each picture show the top score grasp candidate predicted by using corresponding confidence scores generated by force-closure metric or hybrid physical metric. The order of objects is from banana to lion, which is consistent with Table II.

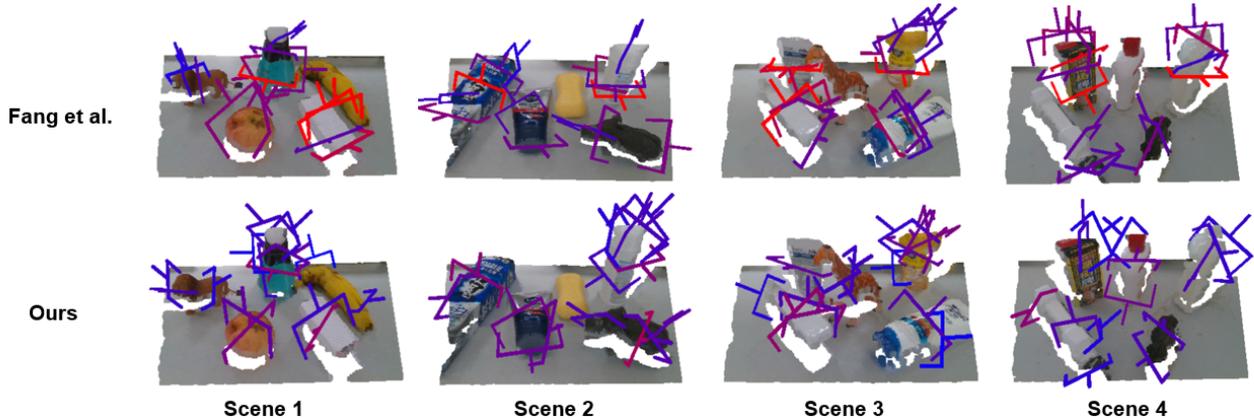


Fig. 7. Scene grasping. There are 15 top grasp candidates after NMS in each picture. The red color represents higher prediction confidence score and the color contrast can only be reflected in the same model. These scenes are consistent with Table III.

another is trained in original score from the force-closure metric. We perform 30 grasp attempts for each object and count the successful grasps. We change the position and pose of the object every time to increase the reliability of the test. The results are reported in Table II. We observe that the introduction of hybrid metrics improves the success rate of grasping in the real world. We show some examples of the compared results under two metrics in the Fig. 6. We can see that the prediction grasps under hybrid physical metric are prone to show the following characteristics: searching the flat contact point, keeping the antipodal direction consistent with the normal direction of the flat contact point, closing to the center of the object, and avoiding too close between gripper end points and contact points.

C. Real Robotic Experiment

The robotic experimental setup includes a UR3 robotic arm and an OnRobot RG2 gripper. The computational resources include NVIDIA TITAN Xp GPU and Intel Xeon E5-2650 CPU. We collect RGB and depth images from an Intel Realsense D435i camera. We set up 4 cluttered table scenes with 5 different objects and conduct several grasp experiments in these scenes. The target of these experiments is to move the objects to the storage box. The camera keeps the same pose in a complete test when collecting the scene images. After predicting grasp results, the top score grasp pose will be given to the robotic arm. It will be regarded as a successful grasp if the object is and placed in the storage box. We set each attempt of one scene for 5 times of image collection and robotic arm motion operation. Therefore, the

single attempt success rate is calculated as the percentage of successful grasps among 5 times. Then we make 10 attempts for each scene and obtain the average success rate. We test [3] and our model respectively. Partial visualization results under four scenes are shown in Fig. 7. The success rate results are reported in Table III. The experiment results prove that our method improves the performance of grasping in the real world.

VI. CONCLUSIONS

We present the hybrid physical metric to evaluate the grasp quality for 6-DoF grasp pose detection. We conduct this metric to generate new grasp confidence scores on the GraspNet-1Billion dataset. We further propose a multi-resolution network FGC-GraspNet to learn these confidence scores better. Altogether, the network evaluation and adequate real robot experiments show that both the hybrid physical metric and the FGC-GraspNet play a positive effect on improving the success rate of grasping. In future work, we aim to apply this framework into complex integrated robotic tasks like feeding a person, cooking or table cleaning.

VII. ACKNOWLEDGMENT

This work was supported by the state key development program in 14th Five-Year under Grant No.2021YFF0602103, 021YFF0602102, 2021QY1702. We also thank for the research fund under Grant No.2019GQG0001 from the Institute for Guo Qiang, Tsinghua University.

REFERENCES

- [1] 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.
- [2] Clemens Eppner, Arsalan Mousavian, and Dieter Fox. Acronym: A large-scale grasp dataset based on simulation. *arXiv preprint arXiv:2011.09584*, 2020.
- [3] Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: A large-scale benchmark for general object grasping. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11444–11453, 2020.
- [4] Corey Goldfeder, Matei Ciocarlie, Hao Dang, and Peter K Allen. The columbia grasp database. In *2009 IEEE international conference on robotics and automation*, pages 1710–1716. IEEE, 2009.
- [5] Minghao Gou, Hao-Shu Fang, Zhanda Zhu, Sheng Xu, Chenxi Wang, and Cewu Lu. Rgb matters: Learning 7-dof grasp poses on monocular rgbd images. *arXiv preprint arXiv:2103.02184*, 2021.
- [6] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *arXiv preprint arXiv:2012.09688*, 2020.
- [7] Kaiyu Hang, Johannes A. Stork, and Danica Kragic. Hierarchical fingertip space for multi-fingered precision grasping. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1641–1648, 2014.
- [8] Yun Jiang, Stephen Moseson, and Ashutosh Saxena. Efficient grasping from rgbd images: Learning using a new rectangle representation. In *2011 IEEE International conference on robotics and automation*, pages 3304–3311. IEEE, 2011.
- [9] Daniel Kappler, Jeannette Bohg, and Stefan Schaal. Leveraging big data for grasp planning. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4304–4311. IEEE, 2015.
- [10] David Kirkpatrick, Bhuvaneshwar Mishra, and Chee-Keng Yap. Quantitative steinitz’s theorems with applications to multifingered grasping. *Discrete & Computational Geometry*, 7(3):295–318, 1992.
- [11] Sulabh Kumra and Christopher Kanan. Robotic grasp detection using deep convolutional neural networks. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 769–776. IEEE, 2017.
- [12] 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.
- [13] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. *Advances in neural information processing systems*, 31:820–830, 2018.
- [14] Hongzhuo Liang, Xiaojian 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.
- [15] 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.
- [16] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof graspnet: Variational grasp generation for object manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2901–2910, 2019.
- [17] Adithyavairavan Murali, Arsalan Mousavian, Clemens Eppner, Chris Paxton, and Dieter Fox. 6-dof grasping for target-driven object manipulation in clutter. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6232–6238. IEEE, 2020.
- [18] Van-Duc Nguyen. Constructing force-closure grasps. *The International Journal of Robotics Research*, 7(3):3–16, 1988.
- [19] Xuran Pan, Zhuofan Xia, Shiji Song, Li Erran Li, and Gao Huang. 3d object detection with pointformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7463–7472, 2021.
- [20] Florian T Pokorny and Danica Kragic. Classical grasp quality evaluation: New algorithms and theory. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3493–3500. IEEE, 2013.
- [21] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas 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*, pages 652–660, 2017.
- [22] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017.
- [23] 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. In *Conference on robot learning*, pages 53–65. PMLR, 2020.
- [24] 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.
- [25] Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017.
- [26] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *arXiv preprint arXiv:1810.04650*, 2018.
- [27] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointcnn: 3d object proposal generation and detection from point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 770–779, 2019.
- [28] Martin Sundermeyer, Arsalan Mousavian, Rudolph Triebel, and Dieter Fox. Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes. *arXiv preprint arXiv:2103.14127*, 2021.
- [29] 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.
- [30] Chao-Ping Tung and Avinash C Kak. Fast construction of force-closure grasps. *IEEE Transactions on Robotics and Automation*, 12(4):615–626, 1996.
- [31] Jacob Varley, Chad DeChant, Adam Richardson, Joaquín Ruales, and Peter Allen. Shape completion enabled robotic grasping. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 2442–2447. IEEE, 2017.
- [32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [33] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9621–9630, 2019.
- [34] Christopher Xie, Yu Xiang, Arsalan Mousavian, and Dieter Fox. The best of both modes: Separately leveraging rgb and depth for unseen object instance segmentation. In *Conference on robot learning*, pages 1369–1378. PMLR, 2020.
- [35] 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 3766–3773. IEEE, 2018.
- [36] Binglei Zhao, Hanbo Zhang, Xuguang Lan, Haoyu Wang, Zhiqiang Tian, and Nanning Zheng. Regnet: Region-based grasp network for end-to-end grasp detection in point clouds. *arXiv preprint arXiv:2002.12647*, 2020.
- [37] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip Torr, and Vladlen Koltun. Point transformer. *arXiv preprint arXiv:2012.09164*, 2020.