# Grasping Detection Network with Uncertainty Estimation for Confidence-Driven Semi-Supervised Domain Adaptation

Haiyue Zhu[1], Yiting Li[2], Fengjun Bai[3], Wenjie Chen[4],
Xiaocong Li[1], Jun Ma[2], Chek Sing Teo[1], Pey Yuen Tao[1], and Wei Lin[1]

*Abstract*— Data-efficient domain adaptation with only a few labelled data is desired for many robotic applications, e.g., in grasping detection, the inference skill learned from a grasping dataset is not universal enough to directly apply on various other daily/industrial applications. This paper presents an approach enabling the easy domain adaptation through a novel grasping detection network with confidence-driven semi-supervised learning, where these two components deeply interact with each other. The proposed grasping detection network specially provides a prediction uncertainty estimation mechanism by leveraging on Feature Pyramid Network (FPN), and the mean-teacher semi-supervised learning utilizes such uncertainty information to emphasizing the consistency loss only for those unlabelled data with high confidence, which we referred it as the confidence-driven mean teacher. This approach largely prevents the student model to learn the incorrect/harmful information from the consistency loss, which speeds up the learning progress and improves the model accuracy. Our results show that the proposed network can achieve high success rate on the Cornell grasping dataset, and for domain adaptation with very limited data, the confidence-driven mean teacher outperforms the original mean teacher and direct training by more than 10% in evaluation loss especially for avoiding the overfitting and model diverging.

## I. INTRODUCTION

Grasping detection is a fundamental problem as most of the sequential robotic manipulation is heavily relying on the successful grasping, e.g. sorting, assembly, etc. The goal for the grasping detection is to find a proper grasp configuration (location and pose), so that the object can be firmly grasped by the gripper. Traditionally, the grasping detection is commonly achieved by human-designed features or 2D/3D model evaluation [1], [2], which calculates the optimal grasping configuration for some certain quality metrics. However, such methods are often tedious and not applicable to novel objects, especially in unstructured environments.

With the recent advances in machine learning, the grasping detection is more and more fulfilled by deep learning [3], [4] to directly learn the grasping candidates from the trials

[1]H. Zhu, X. Li, C. S. Teo, P. Y. Tao, and W. Lin are with the Mechatronics Group, Singapore Institute of Manufacturing Technology (SIMTech), Agency for Science, Technology and Research, Singapore 638075 (e-mail: {zhu_haiyue, li_xiaocong, csteo, pytao, wlin}@simtech.a-star.edu.sg).

[2]Y. Li and J. Ma are with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 117583 (e-mail: yiting_li@u.nus.edu, elemj@nus.edu.sg).

[3]F. Bai is with the Advanced Robotics Application Group, Advanced Remanufacturing and Technology Centre (ARTC), Agency for Science, Technology and Research, Singapore 637143 (e-mail: bai_fengjun@artc.a-star.edu.sg).

[4]W. Chen is with the School of Electrical Engineering and Automation, Anhui University, Hefei, China 230601 (e-mail: wjiechen@yahoo.com.sg).
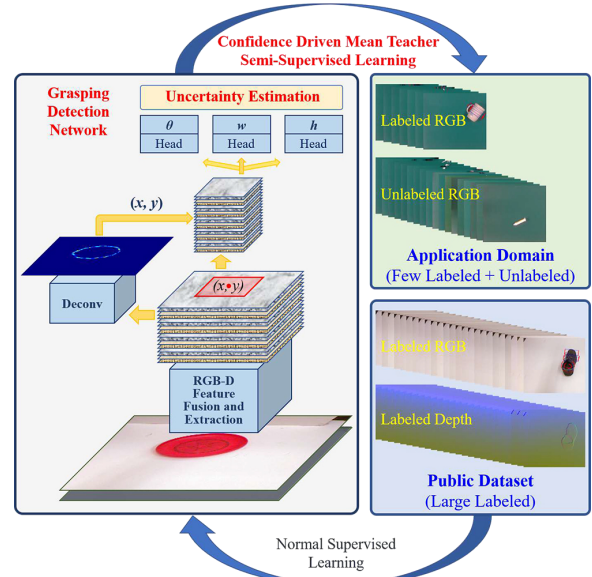
Fig. 1. To facilitate the easy domain adaptation without tedious data labelling, a novel grasping detection network is proposed with a prediction uncertainty estimation mechanism. Furthermore, a confidence-driven mean teacher semi-supervised learning is used to adapt the domain by leveraging the cheap but massive unlabelled data with only a few labelled data.

or human reasoning. The strong generalization capability of deep learning enables such detection networks to tackle well with the novel unseen objects. However, the deep learning for robotic grasping detection is a bit different from the standard computer vision problems as it is quite domain-related and may not be very general. The inference skill acquired from some public datasets may not be directly applicable to other real applications due to the difference in grasping background, sensor, view angle, etc. The traditional solution is to prepare a new domain dataset with full annotations and then apply transfer learning for a direct domain adaptation. However, data preparation and labelling might be quite tedious, especially for the application that may change frequently. As a result, a simple and efficient domain adaptation approach with less labelled data is desired for the robotic grasping application to fill the last-mile gap between artificial intelligence and real-world applications.

In this work, we propose a novel detection architecture that separately treats the whole grasping problem as two subproblems, i.e., location (where) and pose (how), within one network. Different from the widely used region proposal approach [5] in detection networks, the grasping location is

inferred by a heatmap in our network through the encoder-decoder structure, which avoids the disadvantages of the anchor-based approaches such as the imbalance between positive and negative boxes, lots of hyperparameters [6], etc. The location heatmap is generalized from the limited grasping annotations and is able to indicate the grasping feasibility in the continuous space, which facilitates other high-level grasping planning if some additional requirements exist. In parallel, this architecture is also able to infer the optimal grasping pose (angle, width, and height) for every interested grasping location by using the local discriminative feature, thus it constitutes a full grasping detection in an integrated but flexible manner.

Specially, we implement an uncertainty estimation mechanism to evaluate the pose prediction in this architecture. Although Faster/Mask-RCNN also provides a confidence score for every detected box, such a score only represents whether an object is present within the anchor box from the classification view (softmax output), which actually provides no information about the box regression uncertainty. In this work, we propose an uncertainty estimation mechanism based on the Feature Pyramid Network (FPN). In each pyramid stage, the regression of the grasping pose is conducted. If a sample matches the trained distribution well and the prediction is certain, the variance of outputs from different pyramid stages is supposed to be small and vice versa. Moreover, by utilizing the uncertainty estimation mechanism, we propose a confidence-driven mean teacher semi-supervised learning for the grasping domain adaptation. The key idea is that, by leveraging on the cheap unlabelled data through semi-supervised learning, the model consistency is only emphasized for those pseudo-labelled data with high confidence (low uncertainty) instead of them all. This is able to prevent the model from learning the wrong consistency, so that it can improve and speed up the learning process. Overall, the whole picture is illustrated in Fig. 1.

The remainder of this paper is organized as follows. In Section II, we discuss the related work. Section III presents our approaches. The implementation details are introduced in Section IV. The experiments are presented in Section V, and finally, the conclusion is drawn in Section VI.

## II. RELATED WORK

### A. Deep Learning for Grasping Detection

Visual learning approaches for grasping novel objects start in [2], [7]–[9] by sampling and ranking candidate grasps. Deep learning is then introduced for multi-modal information such as RGB-D inputs. Sparse Auto-Encoder (SAE) [10] or stacked SAE [11] with a two-stage approach are used, where a small network is to search a few candidates and a larger network is then to find the top-ranked rectangle. Single-stage methods are proposed [12], [13] with the help of a deep network such as AlexNet [14], ResNet [15], etc., where the grasping detection is formulated as a regression problem to find one optimal location and pose. Multi-grasp detection networks are also proposed for grasping [3], [12], [16], [17], inspired by the recent object detection methods

like YOLO [18] and Faster-RCNN [5], etc. The grasping orientation are also formulated as classification problem by discretization [3], [17]. Grasping detection network for object overlapping scenes are also proposed by finding the Region of Interest (ROI) first and then performing the grasping detection [19].

### B. Semi-Supervised Learning

Semi-supervised learning exploits the unlabelled data to provide the regularization for reducing the model overfitting. It utilizes consistency regularization to make consistent predictions in response to the perturbation of unlabelled samples. Currently, most works are focused on classification. Π-model [20] evaluates the perturbed unlabelled samples twice using stochastic augmentation, dropout and Gaussian noise to minimize their difference and improve the consistency. Mean teacher [21] offers a better teacher model for generating consistency targets through averaging model weights, and its student network is jointly trained with a standard supervised classification loss and an unsupervised consistency loss under different sample augmentations. However, minimizing consistency loss directly might be harmful as the pseudo targets generated by the teacher model are not the real ground truth and might be wrong. Emphasizing the consistency on such noisy data may mislead the training process and degrade the performance, especially when the labelled samples are scarce. In this work, we implement a detection network with prediction uncertainty estimation scheme, our confidence-driven mean teacher model try to leverage the uncertainty metric and only feed the selected confident samples for consistency loss minimization.

### C. Model Uncertainty Estimation

The uncertainty of model prediction is traditionally modelled for the Bayesian neural network by integrating over the posterior distribution over parameters. Recently, it is theoretically proved that Dropout can be used to approximate a model's uncertainty, which can be considered as the Monte Carlo sampling from the posterior distribution of model [22]. It utilizes the variance between multiple predictions with random dropout as a uncertainty metric. In literature, this uncertainty is commonly used for testing purpose only, which is not utilized during model training. In this work, we propose a prediction uncertainty metric based on the consistency between multiple predictions from different pyramid feature stages. Such an uncertainty metric is utilized to rank the detected grasps. More importantly, it is further used to improve the semi-supervised learning for our domain adaptation.

## III. OUR APPROACH

The standard robotic grasping detection problem is formulated [10] as given an RGB-D observation $o$, design a network predictor $f$ to predict the successful grasping rectangle $g$, denoted as $g = f(o)$. The grasp rectangle $g$ is represented with five parameters as

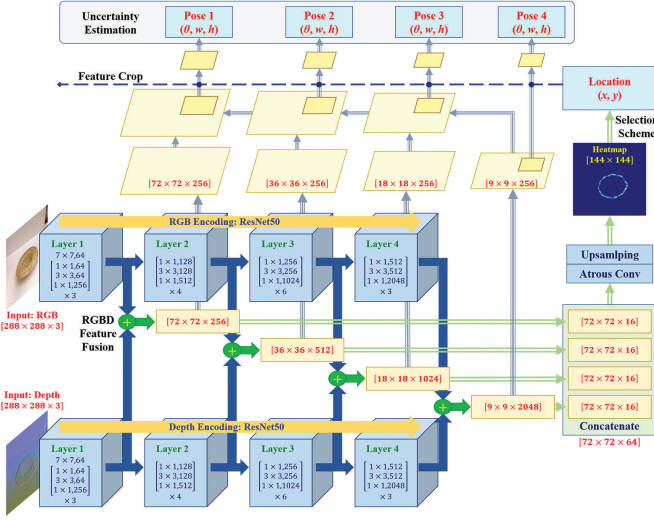$$g = (x, y, \theta, w, h), \tag{1}$$

Fig. 2. Grasping detection network architecture: RGB and depth are encoded in two separate branches by ResNet50 and fused on four different stages, which forms a pyramid feature network. One head decodes the pyramid feature and outputs a grasping location feasibility heatmap. The other four parallel heads crop the local feature around those interested locations on each pyramid feature map and then it predicts the grasping pose separately. The prediction uncertainty is estimated as the variance of four predicted poses.

where $(x, y)$ denotes the center of grasp rectangle, $\theta$ represents its orientation, $w$ denotes the width of gripper plate, and $h$ represents the gripper's opening distance. The grasping detection is a bit special as no real ground truth is fully available because the feasible grasping rectangles might be infinite while the annotations are always limited.

In this work, we propose a grasping detection network to predict the multiple grasping rectangles with the estimated prediction uncertainty, where the network architecture is shown in Fig. 2. Generally speaking, this network treats the grasping detection as two subproblems: "where" are the locations to grasp, denoted as LocNet that

$$(x, y) \leftarrow m = f_l(o \,|\, \phi_s, \phi_l), \qquad (2)$$

and "how" to grasp (referred as pose) for the shortlisted locations, denoted as PoseNet that

$$(\theta, w, h) = f_p(o \,|\, \phi_s, \phi_p, x, y), \qquad (3)$$

instead of directly treating it as a combined "where&how" problem through the overall judgement. Here, $m$ is a grasping location feasibility heatmap, $\phi$ denotes the network parameters, $f_l$ and $f_p$ share the common $\phi_s$ for base feature extraction. Uniquely, the proposed network associates a pose prediction uncertainty metric for every $(\theta, w, h \,|\, o, x, y)$. As a result, the grasping detection network is supposed be more robust by choosing the low-uncertainty candidates to avoid the potential mispredictions.

### A. RGB-D Fusion Based Feature Pyramid Network

To better utilize the RGB and depth information, an RGB-D fusion based on FPN is utilized for feature extraction as

depicted in Fig. 2. The base feature extraction network is commonly shared by both the "where" and "how" detectors $f_l$ and $f_p$, which is relied on ResNet-50 by naturally decomposing it into four block layers $l_k$, $k = 1, 2, 3$, and 4. For both the RGB and depth branches, the feature extraction is separately executed as

$$x_k^r = l_k^r(x_{k-1}^r), \ \ x_k^d = l_k^d(x_{k-1}^d), \qquad (4)$$

where the subscript $r$ and $d$ denote the associations with RGB and depth branches, respectively, and $(x_0^r, x_0^d) = o$ is the input RGB-D pairs. The RGB-D feature fusion is achieved after each ResNet block layer as

$$x_k^f = x_k^r \oplus x_k^d, \qquad (5)$$

where $\oplus$ denotes element-wise summation.

### B. LocNet: Grasping Location Heatmap

The proposed network uses an encoder-decoder Fully Convolutional Network (FCN) head to predict a feasibility heatmap for detecting the feasible/optimal grasping locations, where a higher heat indicates a higher grasping feasibility for this pixel $(x, y)$. For training this FCN head, its loss function is formulated as a binary classification problem for every pixel via the cross-entropy loss. Note that it is not possible to have the real ground truth $m^t$ for all feasible grasping locations, in practical we treat every annotated grasp rectangle's center $(x, y)$ with its surrounding $r$-radius ball as the target $m^t$. Although such a target map is discrete in every annotation, the advantage of using this encoder-decoder structure is that it is a generative detection approach. It can generalize from the sparse grasping location annotations and try to map out all the possible locations in a dense and continuous manner. Such information might also be useful in high-level task planning, e.g., differentiating the handler of a tool to grasp, etc., which is out the scope of this paper. For simplicity, the selection scheme for the obtained heatmap can be filtered by a threshold with the Non-Maximum Suppression (NMS) operation to finalize the shortlisted grasping locations.

### C. PoseNet: Grasping Pose with Uncertainty Estimation

*1) Pyramid Pose Prediction:* With a list of shortlisted locations $(x, y)$ know where to grasp predicted by $f_l$, another head of the proposed network will predict their corresponding pose $(\theta, w, h)$ to further determine how to grasp based on the local discriminative features $x_k^c$. This is a fixed-size sub-feature map around the shortlisted locations $(x, y)$ for every pyramid feature layer $x_k^p$, represented as

$$x_k^c = Cr\Big(x_k^p, (x, y)\Big), \qquad (6)$$

where $Cr$ denotes the crop operation and $k = 1, 2, 3$, and 4. Consequently, four channels of individual fully connected network regression heads $f_{ph}^k$ are utilized to predict the grasping pose on each pyramid feature layer as,

$$\big(\theta_k, w_k, h_k \,|\, x, y\big) = f_{ph}^k(x_k^c) = f_{ph}^k\Big(Cr\Big(x_k^p, (x, y)\Big)\Big), \qquad (7)$$

and finally for a given location $(x, y)$, the mean value is treated as the pose prediction,

$$(\theta, w, h) = mean(\theta_k, w_k, h_k)\big|_{k=1,2,3,4}. \qquad (8)$$

*2) Uncertainty Estimation:* The proposed pyramid pose detection head enables a special mechanism to estimate the pose $(\theta, w, h)$ prediction uncertainty for every $(x, y)$. This is based on the common sense that for a well-trained pose prediction, the prediction results from different pyramid stages should agree with each other, and large variance leads to high uncertainty prediction. As a result, the variance among the predictions from all four pyramid heads is utilized as the uncertainty metric for the pose prediction at $(x, y)$, denoted as $M_{uc}(x, y)$ that,

$$M_{uc}(x, y) = \sum_{\tau = \theta, w, h} var(\tau_k)\big|_{k=1,2,3,4}. \qquad (9)$$

For the shortlisted $(x, y)$ locations, $M_{uc}(x, y)$ helps to identify the most confident predictions of $(\theta, w, h\,|\,x, y)$, which can be ranked as the optimal grasping candidates.

### D. Confidence-Driven Mean Teacher

To achieve domain adaptation via only a small labelled data, the mean teacher semi-supervised learning is adopted in this work. The original mean teacher method is proposed for the classification tasks, we extend its usage on the detection regression network with our new contribution in prediction uncertainty filtering, which we refer it as confidence-driven mean teacher learning. The mean teacher method makes use of the consistency loss to prevent overfitting to limited labelled data. The student learns from the soft pseudo targets provided by the teacher for those unlabelled samples. However, such a pseudo target from the teacher can be wrong itself, the reinforcing on wrong targets will be harmful to the student as it deviates the model convergency. The motivation of confidence-driven mean teacher learning is to let the student only learn those pseudo targets with low uncertainties, which are supposed to be more correct targets. As the training process goes on, the model will become more certain in the new domain, so that more and more pseudo targets with low uncertainties will be added into the training pool gradually for consistency loss optimization. The student model will be benefited from learning those more correct prediction's consistency instead of emphasizing random consistency, thus it speeds up the training progress and improves the overall accuracy.

Let a training set $D_t$ consists of $N_l$ labelled samples and $N_u$ unlabelled samples, denoted as $D_l = \{(e_i, t_i)\}_{i=1}^{N_l}$ and $D_u = \{e_i\}_{i=1}^{N_u}$, respectively, and $D_t = D_l \cup D_u$. Define a teacher model $f^t(\phi')$ and a student model $f^s(\phi)$, these two models share the same network structure but the parameters of teacher model is updated from the student model in every step $k$

$$\phi'_k = \alpha\phi'_{k-1} + (1-\alpha)\phi_k, \qquad (10)$$

where $\alpha$ is a smoothing coefficient hyperparameter. For every training batch $B$, it optimizes the student model by

the minimization the combination of supervision loss and consistency loss

$$\min_{\phi} \left\{ \sum_{e_i \in D_l} \mathcal{L}(f^s(e_i, \phi, \mu), t_i) \right.$$
$$\left. + \sum_{e_i \in \overline{D}_u} \mathcal{L}(f^t(e_i, \phi', \mu') - f^s(e_i, \phi, \mu)) \right\}, \qquad (11)$$

where $\mu$ and $\mu'$ denote the perturbation parameters for the teacher and student models, e.g., augmentation, etc. The uniqueness of our confidence-driven mean teacher exists on

$$\overline{D}_u = \{e_i \,|\, e_i \in D_u \,\&\, M_{uc}(f^t(e_i)) < \overline{T}\}, \qquad (12)$$

which is a subset of $D_u$ that only the samples with low prediction uncertainty $(M_{uc}(f^t(e_i)) < \overline{T})$ will be filtered into $\overline{D}_u$, and $\overline{T}$ is a threshold hyperparameter.

### IV. IMPLEMENTATION DETAILS

### A. Grasping Dataset and Preprocessing

We use the public Cornell grasping dataset (RGB-D) for training the initial grasping detection network. For evaluation of domain adaptation using the proposed confidence-driven mean teacher learning, we also collect a small labelled+unlabelled dataset including 360 images for 18 objects, and only 90 images are labelled with ground truth grasps (five per object). The small new domain dataset is acquired by Intel RealSense D435 RGB-D camera. However, the depth is quite noisy due to the small depth change. Therefore, the small dataset is provided with RGB images only, and the objects, view angle, and background are different from the Cornell grasping dataset for the evaluation of domain adaptation. The original images in both datasets are with 640×480 resolution, they are resized into 456×342 resolution and then cropped into 288×288 as the input to grasping detection network. The depth images are processed into three channels by combining the original channel with Sobel filtered channels in $x$ and $y$.

### B. Training of Grasping Detection Network

The training process of the grasping detection network evolves two steps that the PoseNet and LocNet are trained separately and subsequently. For training the PoseNet, the network is fed with the RGB-D pairs and locations of the grasps in a batch of 16, and then predicts the poses $(\theta, w, h)$ for all pyramid stages. The PoseNet is supervised by minimizing the total losses for all pyramid stages' predictions. The initial learning rate is set as 0.0001 and then gradually decreases, the SmoothL1Loss is used for PoseNet training. Subsequently, the learned coefficients of ResNet-50 layers in PoseNet are passed to LocNet and then keep fixed. The remaining layers of LocNet are further trained by feeding the RGB-D pairs and the targets as described in Section III B, where the CrossEntropyLoss is used and the initial learning rate is 0.0003. The reason we train PoseNet first with ResNet coefficients and then fixed for LocNet is that grasping pose is more important than the location heatmap as we have the certainty based ranking scheme to fine select
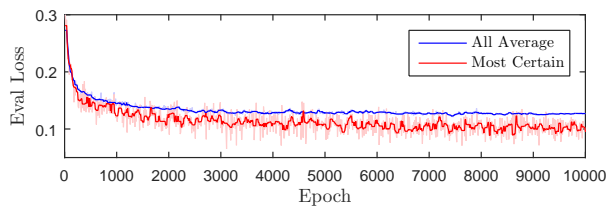
Fig. 3. Comparison between the All Average prediction error against the Most Certain prediction error on the evaluation set.

the optimal location. After this two-step training, the learned coefficients of LocNet and PoseNet are reloaded to the full grasping detection network which can perform the grasping detection tasks with prediction certainty estimation.

### C. Semi-Supervised Domain Adaptation

We apply the proposed confidence-driven mean teacher learning for domain adaptation that enables the grasping detection network learned from the public dataset to perform better for the new application domain with only a small labelled data. The evaluation set contains 54 labelled samples (3 per object) for evaluation. Initially, the LocNet is fine-tuned directly by fully supervised learning using 18 labelled samples. The semi-supervised learning is applied to the head of PoseNet while the ResNet-50 layers keep fixed. For student model PoseNet, the training batch consisting of 6 labelled samples and 2 pseudo labelled samples is fed iteratively to minimize the combined loss, where the perturbation effect is from the sample augmentation. In each step, the teacher model is updated from the new student model, where $\alpha$ is gradually increased from 0.5 to 0.99.

## V. EXPERIMENTS AND DISCUSSIONS

### A. Grasping Detection Network with Certainty Estimation

The evaluation of the proposed grasping detection network is conducted on the Cornell grasping dataset. We first evaluate our proposed uncertain estimation scheme that the lowest-uncertain prediction outperforms the other cases. Fig. 3 shows that the average of the most certain prediction errors against the average of all prediction errors during the whole training progress of PoseNet, where the recorded loss is from the evaluation set that never used for training. It can be seen that the most certain cases are generally better than the average cases around 30% in the loss. Fig. 4 shows the prediction results including location heatmap and detected grasps using the proposed network.

The commonly used grasping evaluation metric is that a candidate grasp is viewed as correct if 1) the difference of angle between predicted grasp $g_p$ and ground truth $g_t$ is within 30°, and 2) the Intersection over Union (IoU) of the predicted grasp $g_p$ and the ground truth $g_t$ is greater than 0.25. Base on this criterion, the proposed grasping detection network achieves 92.2% (177/192) accuracy in the test set (192 samples, 20% image-wise split), where the lowest-uncertain grasp is used for evaluation. However, the above criterion may not able to reflect the real success rate of the
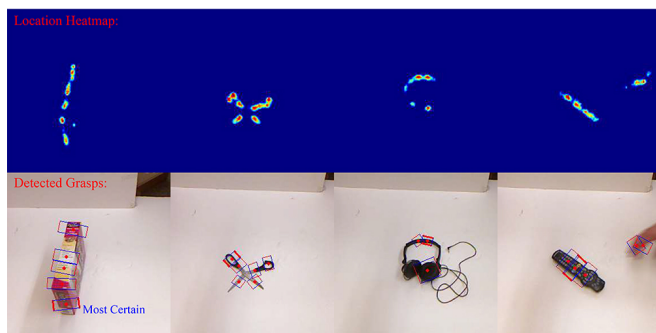


Fig. 4. Predicted location heatmap (above) and grasps (below).



Fig. 5. Evaluation on Cornell grasping dataset: (a) 14 False Negative samples, (b) 1 True Negative sample, and (c) 4 False Positive samples.

grasping as it only accounts for the IoU and angle difference, while the collision between the gripper and object is not considered. By manually analyzing the predictions of all 15 negative samples, we found that 14 cases of them are actually false-negative and only 1 sample is true-negative, shown as in Fig. 5 (a) and (b), respectively. 4 false-positive predictions are found among 177 positive samples, as shown in Fig. 5 (c). Therefore, the real success rate is 187/192=97.4%.

### B. Confidence-Driven Mean Teacher Domain Adaptation

To evaluate the domain adaptation using only limited data, we test the algorithms using 9, 18, and 27 labelled samples for PoseNet, respectively. Fig. 6 plots their comparisons of average loss on a new domain evaluation set. Three training methods are implemented using the same setting, i.e. direct training (fine-tuning) with those labelled data, the original mean teacher to minimize consistency loss for all unlabelled data, and our proposed confidence-driven mean teacher that
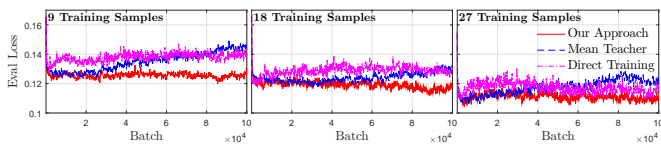
Fig. 6. Comparison of average loss on evaluation set by using (a) 9, (b) 18, and (c) 27 labelled samples for training.

only minimizes consistency loss for these confident samples with low uncertainties. It can be seen from Fig. 6 that the confidence-driven mean teacher can significantly prevent the overfitting in all three cases, while the direct training suffers from the overfitting shortly due to the small number of training sets. Compared with the original mean teacher, the learning speed is faster in the confidence-driven mean teacher, which is because the original mean teacher learns much noisier information with incorrect labels. Moreover, although the original mean teacher performs better than the direct training at the beginning (same level with confidence-driven mean teacher), it can be seen that it disastrously diverges as the training goes on for all three cases, where the reasons are supposed to be same as above due to the emphasizing on the consistency of the wrong pseudo targets.

It is also noted that for real new domain adaptation applications, it is much desired to not have a large labelled evaluation set for having such plots like Fig. 6, all labelled data should be fully used in training to maximize the accuracy. In such a scenario, the confidence-driven mean teacher provides more robust performance in overcoming the overfitting and model diverging. Statistically, by using only 9 labelled data for domain adaptation, the success rate is 145/162=89.5% for adapting from RGB-D to only RGB with different backgrounds, etc. Finally, the grasping detection results using 9 labelled data for domain adaptation are plotted in Fig. 7 with both correct cases (a) and incorrect cases (b), and the network is implemented on the UR robot with Robotiq gripper where the testbed is shown as in Fig. 7(c).

## VI. CONCLUSIONS

This paper addresses the data-efficient domain adaptation for the robotic grasping detection, where only a few labelled data is desired for real applications to minimize the tedious labelling work. We present a grasping detection network with prediction uncertainty estimation and a confidence-driven mean teacher semi-supervised learning algorithm, these two parts are closely interacted to form an integrated solution for data-efficient domain adaptation. Our results show that the proposed detection network is able to perform the grasping detection in high accuracy, and the confidence-driven mean teacher outperforms the original mean teacher and direct training for the detection regression tasks in avoiding the overfitting and model diverging.

## REFERENCES

[1] A. Sahbani, S. El-Khoury, and P. Bidaud, "An overview of 3D object grasp synthesis algorithms," *Robotics and Autonomous Systems*, vol. 60, no. 3, pp. 326–336, 2012.
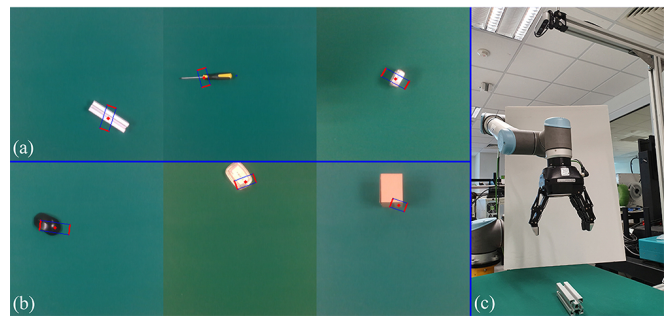


Fig. 7. Prediction evaluation on new domain: (a) correct predictions, (b) incorrect predictions, and (c) the real grasping testbed.

[2] J. Bohg, A. Morales, T. Asfour, and D. Kragic, "Data-driven grasp synthesisa survey," *IEEE Transactions on Robotics*, vol. 30, no. 2, pp. 289–309, Apr 2014.
[3] F. Chu, R. Xu, and P. A. Vela, "Real-world multiobject, multigrasp detection," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3355–3362, Oct 2018.
[4] U. Asif, J. Tang, and S. Harrer, "GraspNet: An efficient convolutional neural network for real-time grasp detection for low-powered devices," in *IJCAI*, 7 2018, pp. 4875–4882.
[5] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *NIPS*, 2015, pp. 91–99.
[6] H. Law and J. Deng, "CornerNet: Detecting objects as paired keypoints," in *ECCV*, Sep 2018.
[7] I. Kamon, T. Flash, and S. Edelman, "Learning to grasp using visual information," in *ICRA*, vol. 3, Apr 1996, pp. 2470–2476.
[8] A. Saxena, J. Driemeyer, and A. Y. Ng, "Robotic grasping of novel objects using vision," *The International Journal of Robotics Research*, vol. 27, no. 2, pp. 157–173, 2008.
[9] Y. Jiang, S. Moseson, and A. Saxena, "Efficient grasping from RGBD images: Learning using a new rectangle representation," in *ICRA*, May 2011, pp. 3304–3311.
[10] 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.
[11] Z. Wang, Z. Li, B. Wang, and H. Liu, "Robot grasp detection using multimodal deep convolutional neural networks," *Advances in Mechanical Engineering*, vol. 8, no. 9, 2016.
[12] J. Redmon and A. Angelova, "Real-time grasp detection using convolutional neural networks," in *ICRA*, May 2015, pp. 1316–1322.
[13] S. Kumra and C. Kanan, "Robotic grasp detection using deep convolutional neural networks," in *IROS*, Sep 2017, pp. 769–776.
[14] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *NIPS*, 2012, pp. 1097–1105.
[15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, Jun 2016.
[16] D. Park, Y. Seo, and S. Y. Chun, "Real-time, highly accurate robotic grasp detection using fully convolutional neural networks with high-resolution images," 2018.
[17] X. Zhou, X. Lan, H. Zhang, Z. Tian, Y. Zhang, and N. Zheng, "Fully convolutional grasp detection network with oriented anchor box," in *IROS*, Oct 2018, pp. 7223–7230.
[18] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *CVPR*, Jun 2016.
[19] H. Zhang, X. Lan, X. Zhou, and N. Zheng, "RoI-based robotic grasp detection in object overlapping scenes using convolutional neural network," *arXiv preprint arXiv: 1808.10313*, 2018.
[20] S. Laine and T. Aila, "Temporal ensembling for semi-supervised learning," *arXiv preprint arXiv: 1610.02242*, 2016.
[21] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *NIPS*, 2017, pp. 1195–1204.
[22] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *ICML*, 2016, pp. 1050–1059.