# Classification based Grasp Detection using Spatial Transformer Network

Dongwon Park[1] and Se Young Chun[1,†]

*Abstract*— **Robotic grasp detection task is still challenging, particularly for novel objects. With the recent advance of deep learning, there have been several works on detecting robotic grasp using neural networks. Typically, regression based grasp detection methods have outperformed classification based detection methods in computation complexity with excellent accuracy. However, classification based robotic grasp detection still seems to have merits such as intermediate step observability and straightforward back propagation routine for end-to-end training. In this work, we propose a novel classification based robotic grasp detection method with multiple-stage spatial transformer networks (STN). Our proposed method was able to achieve state-of-the-art performance in accuracy with real-time computation. Additionally, unlike other regression based grasp detection methods, our proposed method allows partial observation for intermediate results such as grasp location and orientation for a number of grasp configuration candidates.**

## I. INTRODUCTION

Robotic grasping of novel objects is still a challenging problem. It requires to perform robotic grasp detection, trajectory planning and execution. Detecting robotic grasp from imaging sensors is a crucial step for successful grasping. There have been numerous works on robotic grasp detection (or synthesis). In large, grasp synthesis is divided into analytical or empirical (or data-driven) methods [1] for known, familiar or novel objects [2].

Machine learning approaches for robotic grasp detection have utilized data to learn discriminative features for a suitable grasp configuration and to yield excellent performance on generating grasp locations [3], [4]. Since deep learning has been successful in computer vision applications such as image classification [5], [6] and object detection [7], [8], these powerful tools have applied to robotic grasp detection of location and orientation [9], [10].

There are two types of machine learning approaches in robotic grasp detection of location and orientation. One is classification based approach that trains classifiers to discriminate graspable points and orientations for local image patches [3], [4], [9], [11]. In general, sliding window approach generates numerous candidates for robotic grasp configurations with different location, orientation and scale (size of gripper). Then, these image patches are fed into machine learning (or deep learning) classifiers to yield the
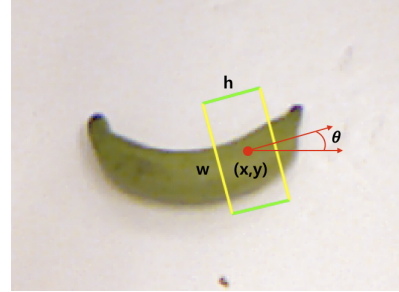


Fig. 1. An example of a robotic grasp detection with five-dimensional grasp representation for a banana. Green lines are two plates of a gripper whose size is $h$, yellow lines are the distance between two plate grippers for grasping, red point is the center location of grasp rectangle $(x, y)$, and red angle $\theta$ is the orientation of the grasp rectangle.

scores of graspability (the higher score, the better graspability). Finally, one candidate image patch with the highest score will be selected. The location, orientation and size of that image patch will be the final grasp detection result. Unfortunately, this approach is in general slow due to brute-force sliding window. Note that classifier based robotic grasp detection is similar to grasp detection methods using 3D grasp simulators [12], [13].

The other is regression based approach that trains a model (parametric or non-parametric, neural network or probability distribution) to yield robotic grasp detection parameters for location and orientation directly [10], [14], [15]. Typical five-dimensional robotic grasp parameters are shown in Fig. 1, where the width $(w)$, height $(h)$, location $(x, y)$ and and orientation $(\theta)$ of a grasp rectangle [11], [9]. Note that with additional depth and surface norm information, these five parameters can be transformed into seven-dimensional robotic grasp representation [11]. Since all grasp parameters are directly estimated from a single image (or a set of multi-modal images), no sliding window is required. Regression based robotic grasp detection methods are usually much faster than classification based methods in term of computation [10]. However, fair comparison between regression and classification based grasp detection does not seem to be well investigated for different computation platforms.

In this paper, we propose a novel classification based robotic grasp detection method using multiple-stage spatial transformer networks (STN). Unlike other black-box regression based grasp detection methods, multiple-stage STN of our proposed method allows partial observation of intermediate grasp results such as grasp location and orientation, for a number of candidates for grasp. We will show that our proposed method achieves state-of-the-art performance in terms of both accuracy and computation complexity. The

[1]Dongwon Park and Se Young Chun are with Department of Electrical Engineering, Ulsan National Institute of Science and Technology (UNIST), Ulsan, Republic of Korea. [†]Email: sychun@unist.ac.kr
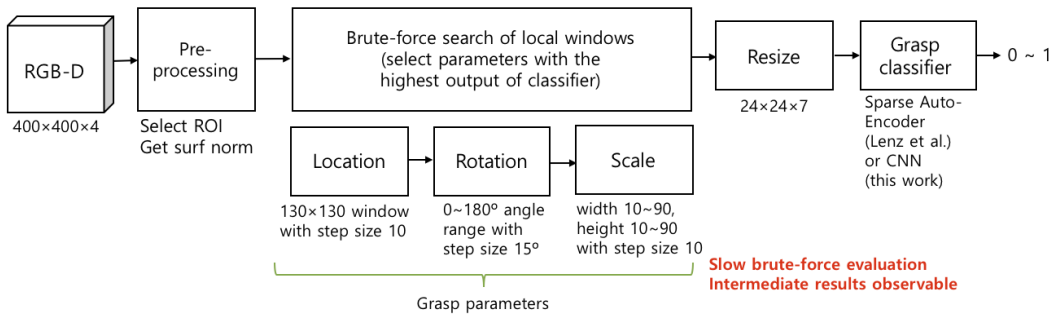
Fig. 2. Classification based robotic grasp detection procedure [9]. Sliding window brute-force search of candidates consumes lots of computation time for different location, orientation, and scale. However, intermediate steps are clearly visible so that good or bad candidates can be observable.

contribution of this paper is as follows:

- A novel multiple-stage STN network is proposed using the original STN [16] and the deep residual network (ResNet) [6] instead of brute force sliding window. Intermediate grasp results are now partially observable in our proposed network.
- A new classification based robotic grasp detection method using our multiple-stage STN is proposed for real-time detection with excellent accuracy. End-to-end training strategy was investigated for our proposed method with high resolution images.
- Extensive comparison of our proposed method with other methods on the same platform was performed with the same size of input image for fair comparison.

## II. RELATED WORK

### A. Deep Learning based Object Detection

There have been much research on object detection using deep learning. Object detection is simply to identify the location and the class of an object (or objects).

There have been several classifier based object detection methods proposed such as region-based convolutional neural network (R-CNN) [17], fast R-CNN [18] and faster R-CNN [7]. The original R-CNN uses a classifier for a local patch of an image with sliding window on the image to identify where and what the object is [17]. Due to time-consuming sliding window with heavy CNN operations, this method is known to be slow. This sliding window approach is similar to the work of Lenz et al. in robotic grasp detection problem [9]. Fast R-CNN significantly reduced computation cost of R-CNN by having sliding window not on an image, but on a feature space [18]. However, due to a large amount of object detection candidates, it was still not a real-time processing. Faster R-CNN proposed region proposal network (RPN) to reduce the amount of object detection candidates so that it significantly improved computation time [7]. RPN generates candidate rectangles of object detection selectively so that faster R-CNN maintains (or improves) detection performance while reduces computation.

Several regressor based object detection methods have also been proposed such as you only look once (YOLO) [19], single shot multibox detector (SSD) [20] and YOLO9000 [8].

Instead of evaluating many object detection candidate windows, these methods process an image only once to estimate object detection rectangles directly so that fast computation is able to be achieved. YOLO used a pre-trained AlexNet [5] to estimate the location and class of multiple objects [19]. SSD further developed regression based object detection to incorporate intermediate CNN features for object detection and improved accuracy and computation speed [20]. Recently, YOLO9000 extended the original YOLO significantly to classify 9000 classes of objects with fast computation and high accuracy [8]. These approaches are similar to the work of Redmon and Angelova [10] in robotic grasp detection.

### B. Deep Learning based Robotic Grasp Detection

Data-driven robotic grasp detection for novel object has been investigated extensively [2]. Before deep learning, there have been some works on grasp location detection using machine learning techniques [3]. Saxena et al. proposed a machine learning method to rank the best graspable location for all candidate image patches from different locations. Jiang et al. proposed to use a five-dimensional robotic grasp representation as also shown in Fig. 1 and further improved a machine learning method to rank the best graspable image patch whose representation includes orientation and gripper distance among all candidates [11].

Since the advent of deep learning [21], robotic grasp detection using deep learning has been investigated for improved accuracy. Lenz et al. proposed a sparse auto-encoder (SAE) to train the network to rank the best graspable image patch with multi-modal information (RGB color, depth, and calculated surface norm) and to apply to robotic grasp detection using sliding window [9]. However, due to time-consuming sliding window process, this method was slow (13.5 sec per image). Moreover, this work was not further extended with simple modification using recent activation functions such as ReLU (Rectifier Linear Unit), which could potentially improve performance significantly [5]. This method can be categorized into classification based grasp detection method. This type of methods allows to observe intermediate steps of grasp detection by showing many candidates with good or bad graspability. Fig. 2 illustrates an example of classification based robotic grasp detection pipeline.

Then, Redmon and Angelova proposed a real-time robotic grasp detection using modified AlexNet [5], [10]. A pre-
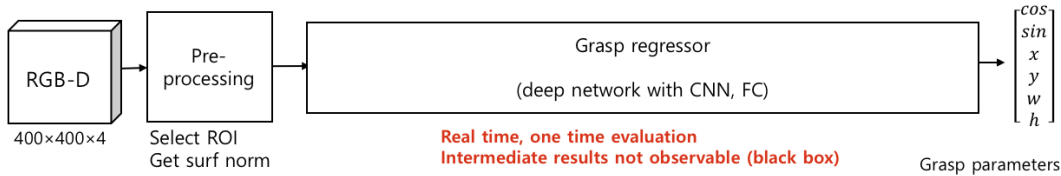
Fig. 3. Regression based robotic grasp detection pipeline [10]. One shot evaluation for the whole image is possible to directly generate grasp configuration such as location, orientation, and scale. However, intermediate steps are not observable.

trained AlexNet was modified to estimate robotic grasp parameters directly for local windows so that no sliding windows is necessary. This approach significantly improved the performance of grasp detection in accuracy and computation time over the work of Lenz *et al.*. To incorporate depth information without changing the network structure much, blue channel was used for depth. No sliding window seems to contribute to the speed up of this method significantly and recent deep network seems to help achieving state-of-the-art accuracy. This method can be categorized into regression based grasp detection method. Unfortunately, this type of regression based methods do not allow to observe intermediate steps of grasp detection. Fig. 3 illustrates an example of regression based robotic grasp detection procedure.

Recently, Asif *et al.* proposed a object recognition and grasp detection using hierarchical cascaded forests using features extracted from deep learning [15]. This work focused on improving accuracy of tasks, rather on improving computation time. Wang *et al.* proposed a real-time classification based grasp detection method using two-stage approach [22]. This method utilized a stacked SAE for classification, which is similar to the work of Lenz *et al.*, but with much more efficient grasp candidate generation. This method utilized several prior information and pre-processing to reduce the search space of grasp candidates such as object recognition result and the graspability of previously evaluated image patches. It also reduced the number of grasp parameters to estimate such as height ($h$) for known gripper and the orientation ($\theta$) that could be analytically calculated from surface norm. This model does not support end-to-end learning for candidate estimation block. Kumra and Kanan proposed a real-time regression based grasp detection method using ResNet for multimodal information (RGB-D). Two pre-trained ResNet-50 networks were used to extract features for RGB color image and for depth image, respectively [14]. Then, a neural network with 3 fully connected layers merged these feature vectors to yield grasp configuration such as width, height, orientation and location. End-to-end training optimized the whole network. However, due to regression based approach, intermediate steps are not observable.

## III. PROPOSED METHOD

### A. Multiple-Stage STN for Robotic Grasp Detection

Instead of slow sliding window for generating many grasp candidates, we propose a multi-stage STN for generating a number of highly selective robotic grasp candidates. This approach seems similar to RPN in [7], but RPN has never used in robotic grasp detection and can not deal with different orientation. STN can explicitly encode spatial transformation including orientation [16]. Our approach is also different from the work of Wang *et al.* [22] that used prior information, while our approach is fully data-driven and can support end-to-end training for fine tuning.

STN consists of localization network to generate transformation parameters, grid generators with the output transformation parameters and sampler to generate warped images or feature maps as shown in Fig. 4 [16]. We modified the original STN by constructing a new localization network using residual blocks [6]. STN can generate all necessary robotic grasp configuration in one network like regression based grasp detection method. However, we adopt multiple-stage approach to generate a number of locations, then to estimate proper orientation for each location, and lastly to fine tune for locations and scale (width and height). Note that this procedure seems similar to human grasp detection. Humans usually find possible grasp locations first, then estimate orientation and scale information. Our multiple-stage STN is illustrated in Fig. 5. Note that since STN is differentiable, it is possible to implement back propagation for our proposed multiple-stage STN.

One of the potential advantages in our multiple-stage STN approach is that intermediate grasp detection steps are partially observable. A number of generated candidates are observable at each stage of our proposed STN. This feature was practically helpful for us to correct for potential errors in grasp detection during our experiment. In addition, it is possible for our proposed multiple-stage STN to be trained end-to-end. This is not only helpful for training robotic grasp detection with ground truth data with labels, but also potentially useful for end-to-end learning from robotic control systems so that the whole network can be improved over trial and error of robots [23].
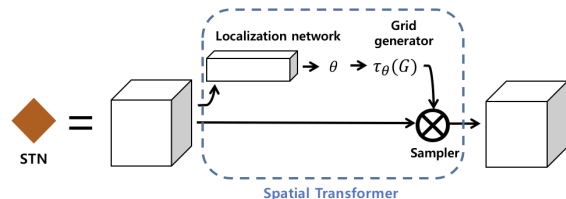


Fig. 4. A typical STN structure proposed in [16]. Localization network generates transformation parameters. Then, STN transforms an input image or feature map using the output parameters. It is possible to back propagate this network since it is differentiable.

Fig. 5. Our proposed classification based robotic grasp detection pipeline. One shot evaluation for the whole image is possible to generate a number of potential grasp candidates so that it is fast as well as intermediate steps are partially observable.

## B. Classification based Grasp Detection using STN

Our proposed grasp detection network using multiple-stage STN is illustrated in Fig. 5. After pre-processing, $STN_{Crop}$ identifies a number of possible grasp locations and generates a set of locations $(x, y)$. Then, for each location, $STN_{rotation}$ is applied to find appropriate orientation $\theta$ for graspable areas. Lastly, $STN_{Scale,Crop}$ determines the width and height for gripper distance and size as well as additional locations $(\delta x, \delta y)$ for fine tuning. Each cropped image patch is fed into grasp classifier to generate the score of graspability. Finally, max pooling select the best graspable configuration among selective candidates.

For training this proposed network, first of all, each component was trained using ground truth data: three STN networks as well as grasp classifier. Then, for end-to-end learning, we focus on training the end block (grasp classifier) first, and then train the second end block ($STN_{Scale,Crop}$) while fixing the previous block, and so on. This fine tuning step was able to improve accuracy further. Note that unlike previous works using $224 \times 224$ input images, our proposed network has relatively high resolution input images with $400 \times 400$, roughly 3 times more pixels than previous cases.

## IV. EXPERIMENTAL RESULTS

Most robotic grasp detection works compared their methods with other previous approaches based on the reported results in papers. While the accuracy comparison may be reasonable, computation complexity comparison may require more careful approach than that due to many crucial factors such as input image size and used GPU spec. In this work, we reproduced all the results of previous methods for comparison on the same platform with the same size data.

### A. Dataset

We trained and tested our proposed classification based robotic grasp detection method using multiple-stage STN with the Cornell grasp detection dataset [9]. This dataset contains 855 images (RGB color and depth) of 240 different objects with the ground truth labels of a few graspable rectangles and a few not-graspable rectangles. To train our

grasp classifiers to assign low graspability score, we additionally generated a number of background image patches with almost all white color. Note that we used cropped images with $400 \times 400$ instead of $224 \times 224$, which is relatively high resolution than the resolution previous works used. It was possible since our proposed method does not require pre-training with massive dataset such as ImageNet.

### B. Training

We split the Cornell grasp dataset into training set and test set (4:1) with image-wise splitting. To train individual blocks of our proposed network as shown in Fig. 4, ground truth labels were transformed into appropriate values for each STN block. We set the output of $STN_{Crop}$ to be 8 so that 4 candidate locations can be generated for potential graspable areas. To generate the initial ground truth for this network, we put the top right ground truth label as the first output of this network and put zero if there are not enough graspable labels. Note that no data augmentation and no pre-training were used.

### C. Evaluation

The same metric for accuracy was used as in [9], [10], [14]. When the difference between the output orientation and



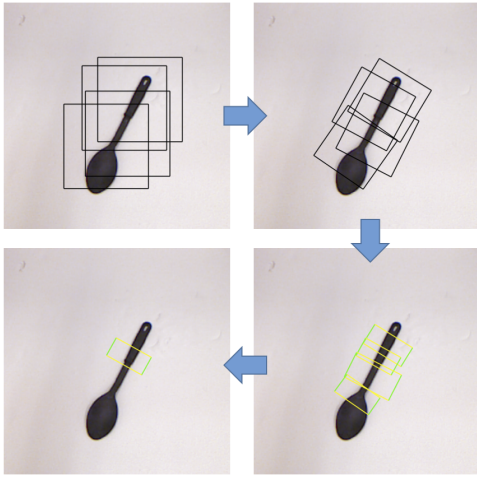Fig. 6. Examples of the Cornell grasp detection dataset [9].

Fig. 7. One result of our proposed grasp detection method. Top left figure is the output of the first STN (4 locations) and top right is the output of the second STN with rotation. Then, the bottom right is the output of the last STN with scaling and fine tuned location. The bottom left is the best graspable configuration among 4 outputs of the last STN.

the ground truth orientation is less than $30^o$, the Jaccard index was measured between the output rectangle and the ground truth rectangle. When the Jaccard index is more than 25%, the output grasp configuration is considered as a good grasp and otherwise bad.

### D. Implementation

The original Lenz *et al.*'s result was reproduced using the MATLAB code provided by the authors [9] (called classification with SAE). We optimized the original code further so that similar computation speed was able to be achieved for larger input images with $400 \times 400$. We also implemented the work of Lenz *et al.* using Tensorflow and replaced SAE with 4-layer CNN (called classification with CNN). For direct regression based grasp detection method that is similar to [10], we implemented it using Tensorflow based ResNet-32 [6]. Unlike the original work using 3 channels (RG and Depth) [10], we implemented the ResNet-32 with 7 channel input so that all multimodal information can be used (RGB color, depth, 3 channel surface norm). Thus, no pre-training was performed for this implementation. Lastly, our proposed methods using single stage STN and multiple stage STN were implemented using Tensorflow. All algorithms were tested on the platform with a single GPU (NVIDIA GeForce GTX 1080 Ti), a single CPU (Intel i7-7700K 4.20GHz) and 32 GB memory.

### E. Qualitative Results of Proposed Method

Fig. 7 shows an example of the step-by-step result using our proposed classification based grasp detection method using multiple-stage STN. The first figure illustrates detected 4 candidate locations from the first STN, then rotated candidates from the second STN for each output of the first STN, scaled and location adjusted candidates from the third STN, and finally chosen one best graspable rectangle from the grasp classifier. Fig. 8 illustrates another example of the result from our proposed method for an object that is difficult
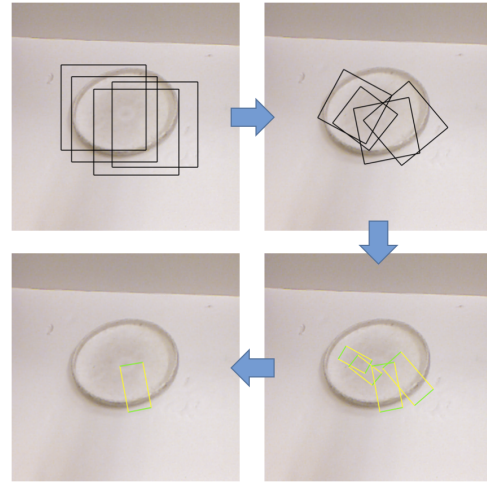


Fig. 8. Another result of our proposed grasp detection method for a difficult object to detect grasp configuration. Top left figure is the output of the first STN (4 locations) and top right is the output of the second STN with rotation. Then, the bottom right is the output of the last STN with scaling and fine tuned location. The bottom left is the best graspable configuration among 4 outputs of the last STN.

to detect robotic grasp configuration with a single output regression model [10]. Thanks to a number of candidates that were generated from the first STN, this case also successfully identify robotic grasp configuration. Note that intermediate steps are observable partially (for 4 candidates in this case) so that it is easier to identify problems for the network than black-box models. Note that this partially observable feature was helpful to design and train our proposed network properly.

### F. Comparison Results for Proposed Method

Fig. 9 illustrates the grasp configuration outputs of classification (SAE), classification (CNN), regression (CNN), and proposed multiple stage STN based grasp detection methods. As also reported in [10], regression based grasp detection method yielded average of all good candidates, while clas-
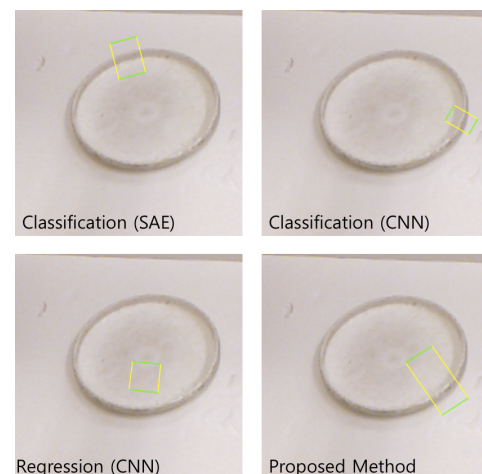


Fig. 9. Four comparison results from classification based grasp detection using sliding window with SAE and CNN, regression based grasp detection and the proposed classification based grasp detection using multi-stage STN.

| Method | Accuracy (%) | Time / Image |
|---|---|---|
| Classification (SAE) | 76.00 | 13 sec |
| Classification (CNN) | 82.53 | 13 sec |
| Regression (CNN) | 70.67 | 11.3 msec |
| Our Single Stage STN | 71.30 | 13.6 msec |
| Our Multiple Stage STN | 89.60 | 23.0 msec |

sification based methods yielded good grasp configurations. Fig. 10 shows that classification (SAE) and regression based methods yielded relatively poor grasp configurations, while our proposed method yielded good grasp configurations. Table I shows that this superior performance of our proposed method is not just for a few images, but for all test images in general. Our method achieved state-of-the-art performance with real-time processing speed.

## V. DISCUSSION AND CONCLUSION

In this paper, we proposed a novel classification based robotic grasp detection method using our multiple stage STN and demonstrated that our proposed method achieved state-of-the-art performance in accuracy and real-time computation time for relatively high resolution images. Our proposed method also has merits such as easy integration with robot control algorithm for end-to-end training, partially observable intermediate steps, and easy training without requiring pre-training with massive amount of data such as ImageNet.

In our results, the accuracy for regression based grasp detection method using ResNet-32 can be further improved when using pre-training. However, in many massive image dataset such as ImageNet, depth images are not available so that it makes challenging to pre-train a model with multimodal information. Regression based grasp detection methods have merit in terms of computation speed, while our proposed method is relatively easy to train without pre-training.



Fig. 10. Four comparison results from classification based grasp detection using sliding window with SAE and CNN, regression based grasp detection and the proposed classification based grasp detection using multi-stage STN.

## 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, Mar. 2012.

[2] 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, Mar. 2014.

[3] 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, Feb. 2008.

[4] J. Bohg and D. Kragic, "Learning grasping points with shape context," *Robotics and Autonomous Systems*, vol. 58, no. 4, pp. 362–377, Apr. 2010.

[5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*, 2012, pp. 1097–1105.

[6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

[7] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems 28*, 2015, pp. 91–99.

[8] J. Redmon and A. Farhadi, "YOLO9000: Better, Faster, Stronger," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6517–6525.

[9] 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, Apr. 2015.

[10] J. Redmon and A. Angelova, "Real-time grasp detection using convolutional neural networks," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2015, pp. 1316–1322.

[11] Y. Jiang, S. Moseson, and A. Saxena, "Efficient grasping from RGBD images: Learning using a new rectangle representation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2011, pp. 3304–3311.

[12] A. T. Miller, S. Knoop, H. I. Christensen, and P. K. Allen, "Automatic grasp planning using shape primitives," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2003, pp. 1824–1829.

[13] B. León, S. Ulbrich, R. Diankov, G. Puche, M. Przybylski, A. Morales, T. Asfour, S. Moisio, J. Bohg, J. Kuffner, and R. Dillmann, "Open-GRASP: A toolkit for robot grasping simulation," in *Simulation, Modeling, and Programming for Autonomous Robots*, 2010, pp. 109–120.

[14] S. Kumra and C. Kanan, "Robotic grasp detection using deep convolutional neural networks," in *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 769–776.

[15] U. Asif, M. Bennamoun, and F. A. Sohel, "RGB-D Object Recognition and Grasp Detection Using Hierarchical Cascaded Forests," *IEEE Transactions on Robotics*, vol. 33, no. 3, pp. 547–564, May 2017.

[16] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial transformer networks," in *Advances in Neural Information Processing Systems 28*, 2015, pp. 2017–2025.

[17] R. Girshick, J. Donahue, and T. Darrell, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 580–587.

[18] R. Girshick, "Fast R-CNN," in *IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1440–1448.

[19] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 779–788.

[20] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single Shot MultiBox Detector," in *European Conference on Computer Vision (ECCV)*, 2016, pp. 21–37.

[21] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.

[22] 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, pp. 1–12, Sept. 2016.

[23] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *Journal of Machine Learning Research*, vol. 17, pp. 1–40, Apr. 2016.