

Learning to Grasp from a Single Demonstration

Pieter Van Molle, Tim Verbelen, Elias De Coninck,
Cedric De Boom, Pieter Simoens, and Bart Dhoedt

Ghent University - imec, IDLab,
Department of Information Technology.
`firstname.lastname@ugent.be`

Abstract. Learning-based approaches for robotic grasping using visual sensors typically require collecting a large size dataset, either manually labeled or by many trial and errors of a robotic manipulator in the real or simulated world. We propose a simpler learning-from-demonstration approach that is able to detect the object to grasp from merely a single demonstration using a convolutional neural network we call GraspNet. In order to increase robustness and decrease the training time even further, we leverage data from previous demonstrations to quickly fine-tune a GraspNet for each new demonstration. We present some preliminary results on a grasping experiment with the Franka Panda cobot for which we can train a GraspNet with only hundreds of train iterations.

Keywords: Learn from demonstration · Deep Learning · Robotics.

1 Introduction

In the advent of Industry 4.0, more and more small and medium enterprises are looking into the adoption of robots to improve their production processes. Increasingly popular are the so-called collaborative robots or cobots. These robots are often lightweight and equipped with force torque sensors, enabling these robots to naturally stop in case of collisions, ensuring safety in human-robot collaboration scenarios [9]. Example cobots that are currently available off-the-shelf are the Kuka LBR series [2], the Universal Robots UR series [3] and the Franka Panda [1], which can be used in a variety of applications such as production line loading and unloading, product assembly, and machine tending [4].

One way to program these collaborative robots is by kinesthetic demonstrations, in which the human operator takes the robot arm and moves it to the desired positions. This reduces the burden of programming the robot, as it is a much more intuitive approach and requires no expert knowledge on robot kinematics or programming code. Although in industry this approach is coined as “learning from demonstration”, it is merely a record and replay feature as opposed to the learning from demonstration research in which generalized policies are trained using machine learning techniques [20]. However, this limits the applicability of currently available systems to cases where positions are fixed relative to the robot. For example, when grasping an object, this object has to

be at the same position for each repetition. Even a small perturbation of the object’s position could harm the system.

In order to mitigate these limitations, one could attach a vision sensor to the robot and use this information to recognize objects, estimate their pose and calculate the best grasp [5]. However, these techniques require lots of supervision, grasp examples and/or training time. In this paper, we propose a grasping approach using neural networks that seamlessly fits in the current established workflow of programming a cobot, and requires only a single demonstration in order to allow perturbations in the grasping position of the target object, up to some extent.

The remainder of this paper is structured as follows. In the next section we give an overview of related work in robotic grasping and learning from demonstration. Next, we propose our approach, which we call GraspNet, in Section 3. We present some preliminary experimental results in Section 4. Finally we discuss our results and conclude with pointers for future work.

2 Related work

Grasping objects with a robotic manipulator is a long standing challenge in research [13]. We focus on data-driven grasping, in which grasping is learned from vision data, either RGB images or depth scans. In particular, we distinguish three types of data-driven grasping: (i) using labeled training data, (ii) using human demonstrations and (iii) using trial and error. For a more in depth survey on data-driven grasping we refer to [5].

The first type assumes a dataset is available with example objects and the corresponding grasp positions, either in the form of 3D meshes [7] or 2D images [19]. Next, a machine learning model such as a neural network can be trained to predict the correct grasp position [18]. In practice this involves carefully collecting and labeling a dataset, on which training can be performed.

A second approach uses human demonstrations [21]. In order to generalize well to situations that are different than the human demonstrations, this is often combined with reinforcement learning techniques [15], or guided policy search [11]. However, for this approach to work, quite a number of demonstrations is usually required in order to generalize well.

Reinforcement learning can also be used in the third type, in which grasping is learned purely from trial and error [16,8]. Another trial and error approach was presented by Levine et al., in which 14 robots collect data over 800,000 grasp attempts [12]. Using this data, a convolutional neural network is trained to predict grasp success, and is then used to implement a controller. A similar approach was presented by Pinto et al. [17]. These trials can also be executed in simulation, after which the learned policies are transferred to the real world [22]. Although these approaches require the least amount of supervision, they are often impractical as they require a lot of trials before the first success.

In this paper we combine elements from the discussed approaches. Similar to [12], we also train a neural network to predict grasp success, which we use

to implement a grasp controller. However, instead of creating a large dataset using trial and error, we collect a single data point per demonstrated grasp. We show that it is possible to train a neural network that predicts the correct grasp position from this single data point through the use of various data augmentation techniques.

3 Learning from a single demonstration

Current collaborative manipulators have a so-called “program by demonstration” feature, with which the operator can easily program a sequence of actions for the robot to execute, by guiding the end effector to the desired positions. This can, for example, be used to program simple pick-and-place tasks. However, at execution time, the robot will merely revisit the programmed positions, without any feedback or closed-loop control about whether the programmed task is actually succeeding. For example, when the object to grasp is not on the exact same position as during the demonstration, this will likely fail.

Our goal is to incorporate visual feedback, by mounting a camera on the end effector of the robot, training a neural network that identifies the object to grasp, and using a closed-loop control algorithm to execute the grasp. In order to mitigate the need of a large-scale dataset for training the neural network, we consider the following assumptions:

- The robot operates in the same workspace as during the demonstrations.
- The object to grasp is the same as during the demonstration.
- The object to grasp can be grasped perpendicular to the workspace plane.
- The object to grasp has a positional offset of at most 8cm, with respect to the original position during demonstration.

Although these assumptions seem limiting, they still cover most use cases in an industrial environment, where objects need to be picked from and placed in well

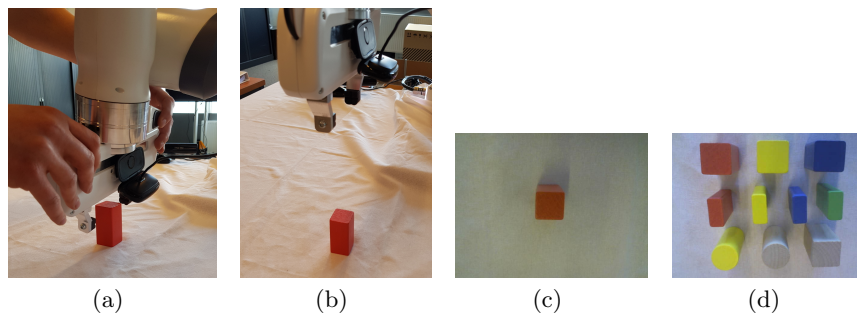


Fig. 1. First, the operator guides the robot to the correct grasp pose (a). Next, the robot moves up (b) and the top-down camera captures a single frame used for training (c). We consider 10 types of blocks to grasp (d).

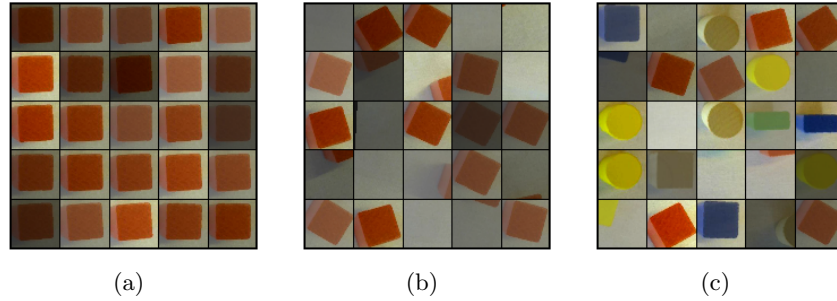


Fig. 2. From the demonstration camera frame, we generate random positive (a) and negative (b) samples by cropping, rotating and adapting brightness and contrast. By including the demonstrations for the other blocks, we can generate additional negative samples (c).

defined bin areas, but are not necessarily nicely aligned on fixed positions within these bins.

3.1 Setup

Our setup is shown in Figure 1 and consists of the Franka EMIKA Panda cobot, with a camera mounted on its end effector. To record a demonstration, the operator simply guides the gripper to the preferred grasp pose (a). Next, the robot will hover this position on a fixed height (b), and one camera frame is recorded (c). On this camera frame, the object to grasp will be positioned in the center of the image. As target objects, we currently use toy blocks of different shapes and colors. We gather demonstrations for the 10 types of blocks shown in (d).

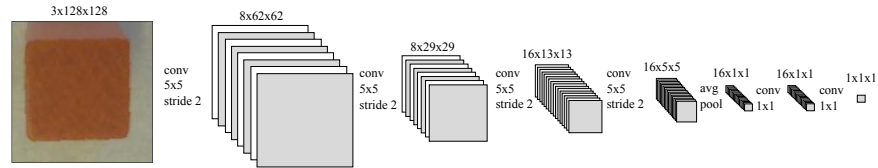


Fig. 3. GraspNet architecture: an 128×128 image crop is processed by 4 convolutional layers with 8, 8, 16 and 16 filters of size 5×5 and stride 2, followed by an average pooling layer and two fully connected layer with 16 hidden units. The fully connected layers are implemented as 1×1 convolutions and the output is a sigmoid neuron estimating grasp success.

3.2 GraspNet

We train a convolutional neural network, dubbed *GraspNet*, to detect the correct grasp pose for a target object based on a camera frame, requiring but a single human demonstration. A well-known technique for extending a dataset and improving neural network training is data augmentation [6], where images are perturbed to generate additional examples from the same underlying class. In this work, we rely on extreme data augmentation in order to create a “very large” dataset, starting from the single demonstration camera frame.

From the demonstration camera frame, with a resolution of 640×480 pixels, we generate a train set by taking random crops of 128×128 pixels. Positive samples are generated by taking center crops with a small random rotation (sampled uniformly between -3 and 3 degrees), while any other crop and rotation is used as negative sample. All samples are further randomized by applying random perturbations on brightness and contrast, both sampled uniformly between 0.5 and 1.5. Examples of positive and negative samples are shown in Figure 2.

Our neural network architecture is depicted in Figure 3. It consists of four convolutional layers, of which the final layer is pooled using average pooling. We include two fully connected layers at the end of the network, implemented as 1×1 convolutions. All hidden layers have rectified linear units (ReLU), and the network ends with a single sigmoid neuron to classify each crop as positive or negative.

3.3 Grasp controller

By implementing the fully connected layers as 1×1 convolution kernels, we can easily apply this to an image with arbitrary size. The output is then a two-dimensional feature plane, which we interpret as an activation map that indicates the presence and position of the target grasp object, as seen in Figure 4.

We implement a Cartesian velocity controller for the robot arm, that calculates the direction vector v from the center of the image to the point of highest activation, and maps this to movement in the workspace plane. Once the highest

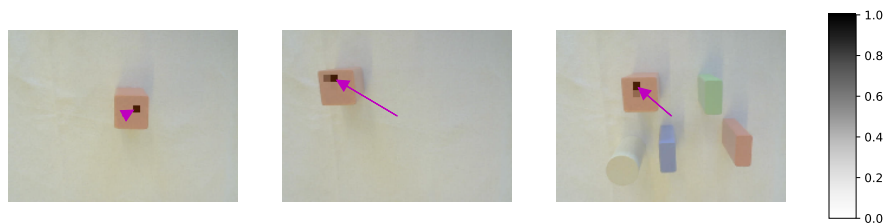


Fig. 4. Examples of camera images with activation map overlays, generated by running the images through GraspNet, and upscaling the output feature planes to match the input dimensions. Darker regions correspond to higher activations. The direction vector v is represented by the magenta arrow. The big red block is always the target object.

activation is in the center of the image, the arm moves straight down and closes the gripper.

In order to also incorporate rotational information, we generate a batch of rotations of the camera image (e.g. ranging from -50 degrees to 50 degrees, with a step size of 10), and forward this batch through GraspNet, resulting in a batch of activation maps. Next, we rotate in the direction of the rotation with the highest maximum activation. This process is repeated until the unrotated image has the highest maximum activation.

By constantly streaming the camera images through GraspNet and acting accordingly we get a closed-loop controller that can successfully grasp the object once it is in the field of view of the camera.

3.4 Grasping n objects

We can easily extend our approach when learning to grasp multiple object types, by providing a demonstration and training a separate GraspNet for each object type. However, when two objects look similar (for example the two red or blue blocks in our experiments), these might be difficult to distinguish, as we only train GraspNet from a single demonstration sample. In order to make GraspNet more robust, we combine all demonstration frames: when training for a certain target object, we extend our train set by regarding positive samples for other objects as negative samples for the target object. Examples of such negative samples are shown in Figure 2(c).

The main drawback of our approach is that we have to train a neural network from scratch for each demonstration. In order to mitigate this, we further improve our methodology by also learning a set of initialization parameters that can be fine-tuned quickly for a new demonstration. We apply the Reptile algorithm [14], a recently proposed approach for meta-learning for few-shot classification. Suppose we already have a set of n demonstrations, we can train a set of initialization parameters ϕ by iteratively sampling one of the demonstrations, training GraspNet weights W for this demonstration by applying k gradient descent steps, and then updating $\phi \leftarrow \phi + \epsilon(W - \phi)$, where ϵ is the outer step size. This means we push ϕ to an optimum in which for all n tasks we can get a high performing set of weights W after k additional gradient descent steps.

4 Experiments

We apply the dataset generation schemes described in Section 3 to generate a basic train set (using the single demonstration frame) and an extended train set (using all demonstration frames) for each of the 10 target blocks. Using these train sets, we train two GraspNet configurations for each block: “single”, in which the basic train set for the target block is used, and “multi”, where the extended train set is used. Each GraspNet is trained for 512 iterations, with batches of 64 random samples. Parameter updates are done using the Adam algorithm [10], with a step size of 0.001. We also implement the Reptile algorithm, for which

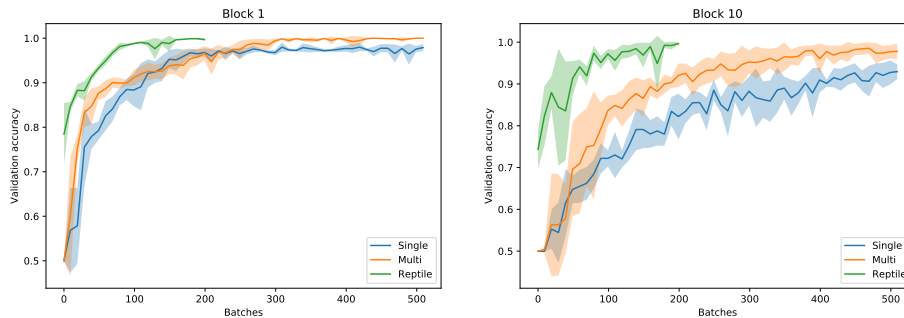


Fig. 5. Training performance for each configuration on an easy task (block 1) and a more difficult task (block 10). Averages over 5 random seeds are shown.

we first train initialization parameters on nine of the demonstrations, and then evaluate for the remaining block. In this case, we first train a GraspNet for 250 outer iterations, with an initial outer step size of 0.6, which is linearly annealed to 0. We use meta-batches, where we train separately on four tasks, and average the update directions. Each task is trained for 10 iterations with a batch size of 10. Parameters are updated using the Adam algorithm, with a step size of 0.001, and $\beta_1 = 0$, to disable momentum, as done by [14]. Next, we train for 200 iterations of batch size 64 on the target task. All training is executed on an Nvidia Titan X GPU. The classification accuracy on a hold-out validation set (comprised of 200 random positive and 200 random negative samples, generated using the extended dataset generation scheme) for two blocks is plotted in Figure 5.

The results show that training using data of all demonstrations (multi) yields indeed better results, especially for the blocks that are more difficult to distinguish. We also see that training with Reptile indeed converges much faster, and in some cases even yields the best results.

We report the final classification accuracies for each of the configurations in Table 1. For each target block, we use a test set of 5,000 random positive and 5,000 random negative samples, again, generated using the extend dataset generation scheme. We see again that multi performs better or equal compared to single, but also that Reptile performs on par (or better), although this only requires a fraction of the training iterations for fine-tuning.

5 Discussion

We also conduct some initial experiments using the real robot controlled by an Intel NUC equipped with an Intel i5-5250U CPU. When we do not take rotation into account, the robot is able to robustly pick up the block using GraspNet. When applying random rotations, we generate a batch of nine planes where we rotate the camera frame in steps of 10 degrees, and feed it through GraspNet. In this case, our robot successfully detects the correct grasp orientation, but the

Table 1. Test set accuracies per configuration, for each target block, averaged over 5 random seeds.

Block	Single	Multi	Reptile
1	0.989 ± 0.004	0.999 ± 0.001	0.997 ± 0.001
2	0.985 ± 0.009	0.991 ± 0.015	0.989 ± 0.007
3	1.000 ± 0.000	0.995 ± 0.010	0.958 ± 0.007
4	0.976 ± 0.015	0.993 ± 0.009	0.999 ± 0.000
5	0.969 ± 0.021	0.999 ± 0.001	0.983 ± 0.003
6	0.995 ± 0.008	0.959 ± 0.072	0.954 ± 0.046
7	0.910 ± 0.055	0.994 ± 0.007	0.938 ± 0.077
8	0.957 ± 0.007	0.997 ± 0.002	0.974 ± 0.024
9	0.887 ± 0.035	0.960 ± 0.021	0.925 ± 0.068
10	0.927 ± 0.018	0.974 ± 0.022	0.993 ± 0.002

closed-loop controller fails to move the gripper to the exact grasp location. This is due to the fact that processing the batched input takes up to 700ms to process which is too high for accurate control. We are now looking into extending our setup using a separate, GPU-enabled machine for processing the camera frames.

Besides these promising results, we are aware of the current limitations of our approach. For example, as we only use a few data examples, this reduces the applicability to new and unseen situations. Of course this can be mitigated by adding additional data samples, or by using available grasp datasets [19]. This would be especially appealing for usage of the algorithm, in which we can train once on a large dataset to obtain suitable initialization parameters to quickly fine-tune a GraspNet for each new demonstration.

One of the biggest advantages of our approach is that we can integrate it seamlessly with the currently common programming-by-demonstration workflow. A non-technical operator can still easily program the robot, while alleviating the need of supplying objects to grasp on the exact same position.

6 Conclusion

In this paper, we have proposed a technique for learning to grasp an object from a single demonstration from a vision sensor. For each demonstration, we trained a GraspNet, a convolutional neural network that predicts grasp success based on an image patch. Using such a GraspNet, we developed a closed-loop robot controller to grasp the object whenever it is in view of the camera. We show that we can in fact train a GraspNet from a single image, and that we can increase robustness by also including data from a few other demonstrations. Moreover,

using the Reptile meta-learning algorithm, we show that we can quickly fine-tune a GraspNet for a new demonstration.

In future work, we plan to extend our approach to grasp more common objects instead of merely toy blocks, and integrate GraspNet in a programming-by-demonstration platform for a real collaborative industrial workspace. This will allow to test our approach on production-grade hardware platforms.

Acknowledgements

Cedric De Boom is funded by a PhD grant of the Research Foundation - Flanders (FWO). We would like to thank Nvidia Corporation for their generous donation of a Titan X GPU.

References

1. Franka EMIKA. <http://www.franka.de/>, last accessed April 6, 2018
2. Kuka AG. <http://www.kuka.com/>, last accessed April 6, 2018
3. Universal robots. <http://www.universal-robots.com/>, last accessed April 6, 2018
4. Bloss, R.: Collaborative robots are rapidly providing major improvements in productivity, safety, programming ease, portability and cost while addressing many new applications. *Industrial Robot: the international journal of robotics research and application* **43**(5), 463–468 (2016)
5. Bohg, J., Morales, A., Asfour, T., Kragic, D.: Data-driven grasp synthesis - a survey. *IEEE Transactions on Robotics* **30**(2), 289–309 (April 2014)
6. Chatfield, K., Simonyan, K., Vedaldi, A., Zisserman, A.: Return of the devil in the details: Delving deep into convolutional nets. In: *BMVC* (2014)
7. Goldfeder, C., Ciocarlie, M., Dang, H., Allen, P.K.: The columbia grasp database. In: *Proceedings of the 2009 IEEE International Conference on Robotics and Automation*. pp. 3343–3349 (2009)
8. Gu, S., Holly, E., Lillicrap, T., Levine, S.: Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*. pp. 3389–3396 (May 2017)
9. Haddadin, S., Albu-Schaffer, A., Luca, A.D., Hirzinger, G.: Collision detection and reaction: A contribution to safe physical human-robot interaction. In: *2008 IEEE/RSJ International Conference on Intelligent Robots and Systems*. pp. 3356–3363 (Sept 2008)
10. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
11. Levine, S., Finn, C., Darrell, T., Abbeel, P.: End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research* **17**(1), 1334–1373 (Jan 2016)
12. Levine, S., Pastor, P., Krizhevsky, A., Ibarz, J., Quillen, D.: Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International Journal of Robotics Research* (2017)
13. Nguyen, V.D.: Constructing force-closure grasps. In: *Proceedings. 1986 IEEE International Conference on Robotics and Automation*. vol. 3, pp. 1368–1373 (Apr 1986)
14. Nichol, A., Schulman, J.: Reptile: a scalable metalearning algorithm (March 2018)

15. Pastor, P., Hoffmann, H., Asfour, T., Schaal, S.: Learning and generalization of motor skills by learning from demonstration. In: 2009 IEEE International Conference on Robotics and Automation. pp. 763–768 (May 2009)
16. Peters, J., Schaal, S.: Reinforcement learning of motor skills with policy gradients. *Neural Networks* **21**(4), 682–697 (2008)
17. Pinto, L., Gupta, A.: Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In: 2016 IEEE International Conference on Robotics and Automation (ICRA). pp. 3406–3413 (May 2016)
18. Redmon, J., Angelova, A.: Real-time grasp detection using convolutional neural networks (December 2014)
19. Saxena, A., Driemeyer, J., Ng, A.Y.: Robotic grasping of novel objects using vision. *The International Journal of Robotics Research* **27**(2), 157–173 (2008)
20. Schaal, S.: Learning from demonstration. In: Proceedings of the 9th International Conference on Neural Information Processing Systems. pp. 1040–1046 (1996)
21. Tegin, J., Ekvall, S., Kragic, D., Wikander, J., Iliev, B.: Demonstration-based learning and control for automatic grasping. *Intelligent Service Robotics* **2**(1), 23–30 (Jan 2009)
22. Tobin, J., Fong, R., Ray, A.K., Schneider, J., Zaremba, W., Abbeel, P.: Domain randomization for transferring deep neural networks from simulation to the real world. 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) pp. 23–30 (2017)