

---

# Training Dense Object Nets: A Novel Approach

---

**Anonymous Author(s)**

Affiliation  
Address  
email

## Abstract

1 We present a novel framework for mining dense visual object descriptors produced  
2 by Dense Object Nets (DON) without explicitly training DON. DON’s dense visual  
3 object descriptors are robust to changes in viewpoint and configuration. However,  
4 training DON requires image pairs with correspondence mapping, which can be  
5 computationally expensive and limit the dimensionality and robustness of the de-  
6 scriptors, limiting object generalization. To overcome this, we propose a synthetic  
7 augmentation data generation procedure and a novel deep learning architecture  
8 that produces denser visual descriptors while consuming fewer computational re-  
9 sources. Furthermore, our framework does not require image-pair correspondence  
10 mapping and demonstrates its one of the applications as a robot-grasping pipeline.  
11 Experiments show that our approach produces descriptors as robust as DON.

12 **1 Introduction**

13 Creating a general-purpose robot capable of carrying out practical activities like Chappie [1] or  
14 C-3PO [2], is a long-standing objective of robotics and robotic manipulation. While advancements  
15 have been made recently in adjacent domains, achieving this goal remains a work in progress. For  
16 instance, AlphaGo [3], a game-playing artificial intelligence system trained entirely on self-play,  
17 defeated the world’s best human Go player at the time. Subsequently, Silver et al. [4] developed  
18 artificial intelligence algorithms that mastered the game of chess, Go, World of Warcraft [5], and  
19 Shogi, surpassing human playing expertise. Most of these algorithms learn directly from visual  
20 data such as gameplay recordings or online video streams, emphasizing the importance of visual  
21 data in AI. Meanwhile, the launch of AlexNet [6] in 2012 transformed the field of computer vision.  
22 Other visual tasks, such as semantic segmentation [7], object identification and recognition [8],  
23 and human pose estimation [9], have also witnessed significant gains in recent years. In robotics,  
24 significant breakthroughs have been made, ranging from self-driving cars to humanoid robots capable  
25 of performing complex tasks using cameras and other vision sensors. Despite these advancements,  
26 the most frequently used robotic manipulation systems have not evolved much in the previous 30  
27 years. Typical auto-factory robots continue to perform repetitive operations such as welding and  
28 painting, with the robot following a pre-programmed course with no feedback from the surroundings.  
29 If we want to increase the utility of our robots, we must move away from highly controlled settings  
30 and robots that perform repetitive actions with little feedback or adaptability capabilities. Liberating  
31 ourselves from these constraints of controlled settings-based manufacturing would allow us to enter  
32 new markets, as witnessed by the proliferation of firms [10] competing in the logistics domain.

33 The ideal object representation for robot grasping and manipulation tasks remains to be engineered  
34 today. Existing representations may not be suitable for complex tasks due to limited capabilities  
35 of understanding an object’s geometrical and structural information. In 2018, Florence et al. [11]  
36 introduced a novel visual object representation to the robotics community, terming it “dense visual  
37 object descriptors”. DON, an artificial intelligence framework proposed by Florence et al. [11]  
38 produces dense visual object descriptors. In detail, the DON converts every pixel in the image

39 ( $I[u, v] \in \mathbb{R}^3$ ) to a higher dimensional embedding ( $I_D[u, v] \in \mathbb{R}^D$ ) such that  $D \in \mathbb{N}^+$  consuming  
 40 image-pair correspondences as input yielding pixelwise embeddings which are nothing but dense  
 41 local descriptors. The dense visual object descriptor generalizes an object up to a certain extent  
 42 and has been recently applied to rope manipulation [12], block manipulation [13], robot control  
 43 [14], fabric manipulation [15] and robot grasp pose estimation [16, 17]. Adrian et al. [17] further  
 44 demonstrated that DON can be trained on synthetic data and still generalize to real-world objects.  
 45 Furthermore, Adrian et al. [17] demonstrated that the quality of descriptors produced by the DON  
 46 framework depends on the higher or longer embedding dimension. We tried training the DON on  
 47 a computation device equipped with NVIDIA RTX A6000 GPU with 48GB VRAM. However, we  
 48 could not train the DON to produce a higher embedding dimension due to the limited VRAM. The  
 49 DON framework is computationally expensive, as shown in Table 1, and limits the user to generalize  
 50 objects to a certain extent making it difficult to use as a robot grasping pipeline in real-world logistics  
 51 and warehouse automation scenarios.

Table 1: Benchmark of DON framework trained on GPU with 48GB VRAM with 128 image-pair correspondences, batch size of 1 and “Pixelwise NTXENT Loss” [17] as a loss function.

GPU VRAM consumption to train DON				
Descriptor Dimension	3	8	16	32
VRAM Usage (GB)	9.377	13.717	20.479	30.067

52 To overcome the computation resource limitation to produce denser visual object descriptors, we  
 53 propose a novel framework to train and extract dense visual object descriptors produced by DON,  
 54 which is computationally efficient.

## 55 2 Related Work

56 We are solely interested in computing dense visual object descriptors of an object. The DON training  
 57 strategy in [11] relies on the depth information for computing correspondences in an image pair  
 58 using camera intrinsics and pose information [18]. However, when employing consumer-grade depth  
 59 cameras for capturing the depth information, the depth cameras capture noisy depth in cases of tiny,  
 60 reflecting objects, which are common in industrial environments. In the meantime, Kupcsik et al. [16]  
 61 used Laplacian Eigenmaps [19] to embed a 3D object model into an optimally generated embedding  
 62 space acting as an target to train DON in a supervised fashion. The optimal embeddings brings  
 63 in more domain knowledge by associating 3D object model to images views. Kupcsik et al. [16]  
 64 efficiently apply it to smaller, texture-less and reflective objects by eliminating the need of the depth  
 65 information. Kupcsik et al. [16] further compare training strategies for producing 6D grasps for  
 66 industrial objects and show that a unique supervised training approach increases pick-and-place  
 67 resilience in industry-relevant tasks.

68 Florence [20] has found that the pixelwise contrastive loss function used to train DON might not  
 69 perform well if a computed correspondence is spatially inconsistent (analogously to the case of noisy  
 70 depth information). This further highlights that the precision of contrastive-trained models can be  
 71 sensitive to the relative weighting between positive-negative sampled pixels. Instead, the Florence [20]  
 72 introduces a new continuous sampling-based loss function called “Pixelwise Distribution Loss”. The  
 73 pixelwise distribution loss is much more effective as it is a smooth continuous pixel space sampling  
 74 method compared to the discrete pixel space sampling method based on pixelwise contrastive  
 75 loss. The pixelwise distribution loss regresses a set of probability distribution heatmaps aiming to  
 76 minimize the divergence between the predicted heatmap and the ground truth heatmap mitigating  
 77 errors in correspondences. Futhermore, the pixelwise distribution loss does not need non-matching  
 78 correspondences compared to the the pixelwise contrastive loss. Differently, Hadjivelichkov and  
 79 Kanoulas [21] extends the DON training using semantic correspondences between objects in multi-  
 80 object or cluttered scenes overcoming the limitations of [18, 19]. The authors, Hadjivelichkov and  
 81 Kanoulas [21] employ offline unsupervised clustering based on confidence in object similarities  
 82 to generate hard and soft correspondence labels. The computed hard and soft labels lead DON in  
 83 learning class-aware dense object descriptors, introducing hard and soft margin constraints in the  
 84 proposed pixelwise contrastive loss to train DON. Further eliminating the need for camera pose  
 85 and intrinsic information along with depth information to compute correspondences in an image

pair, Yen-Chen et al. [22] used NeRF [23] to train DON. The NeRF [23] recreates a 3D scene from a sequence of images captured by the smartphone camera. The correspondences are extracted from the synthetically reconstructed scene to train DON. Recently, based on SIMCLR inspired frameworks [24, 25], Adrian et al. [17] introduced similar architecture and another novel loss function called “Pixelwise NT-Xent loss” to train DON more robustly. The pixelwise ntxent loss consumes synthetic correspondences independent of depth cameras computed from image augmentations to train DON. Adrian et al.’s experiments show that the novel loss function is invariant with respect to the batch size. Additionally adopted “*PCK@k*” metric has been adopted as in preceedings [26, 27] to evaluate and benchmark DON on cluttered scenes previously not benchmarked.

In the proposed framework we do not use any loss functions in [11, 20, 16, 17, 21, 22] to train DON however we adopt the network architecture from [11] and train on the task of the “KeypointNet”[28] with adaption of the loss functions proposed in [28, 29].

### 3 Methodology

#### 3.1 Dataset Engineering

We have chosen the cap object for creating synthetic dataset as the cap mesh models are readily available in the “Shapenet” library [30] as it contains rich object information including textures. Furthermore, we choose 5 cap models from the Shapenet library and use Blenderproc [31] to generate the synthetic dataset. For each cap model we save one RGB image, mask and depth in the synthetic scene. Additionally, we employ synthetic augmentations as proposed in [17] to synthetically spatial augment cap’s position and rotation in an image including background randomization using Torchvision [32] library. To generate camera poses for different viewpoints, an augmented image-pair is sampled randomly and image-pair correspondences is computed <sup>1</sup>as illustrated in the Figure 1. Using depth information we project the computed correspondences to camera frame and compute relative transformation between two camera-frame coordinates of the correspondences using Kabsch’s transformation [33].



Figure 1: Depiction of image synthetic spatial augmentation and correspondences mapping in an image-pair. The colored encoded dots in the figure represents correspondences in an image-pair.

#### 3.2 Framework & Mining Strategy

As a backbone, we employ ResNet-34 architecture [34]. We preserve the last convolution layer and remove the pooling and linear layers. The backbone downsamples the RGB image  $I_{RGB} \in \mathbb{R}^{H \times W \times 3}$  to dense features  $I_d \in \mathbb{R}^{h \times w \times D}$  such that  $h \ll H, w \ll W$  and  $D \in \mathbb{N}^+$ . We upsample the dense features from the identity layer (being identity to the last convolution layer in the backbone) as illustrated in the Figure 2 in page 4 as follows:

$$f_U : I \in \mathbb{R}^{h \times w \times D} \rightarrow I_D \in \mathbb{R}^{H \times W \times D}. \quad (1)$$

The upsampled dense features is extracted and treated as dense visual local descriptors produced from the DON in otherwords we extract or mine the representations from the backbone. Similarly as in [28], we stack spatial-probability regressing layer and depth regressing layer on top of the identity layer to predict  $N \in \mathbb{N}^+$  number of keypoint’s spatial-probability as follows:

$$f_S : I_d \in \mathbb{R}^{h \times w \times D} \rightarrow I_s^N \in \mathbb{R}^{h \times w \times N}, \quad (2)$$

<sup>1</sup>GitHub Link:

<https://github.com/KanishkNavale/Mapping-Synthetic-Correspondences-in-an-Image-Pair>

121 and depth as follows:

$$f_D : I_d \in \mathbb{R}^{h \times w \times D} \rightarrow I_{\hat{d}} \in \mathbb{R}^{h \times w \times N}. \quad (3)$$

122 We incorporate continuous sampling method  $f_E$  from [20, 28] to convert the upsampled predicted  
123 spatial-probability and depth of a keypoint to spatial-depth expectation as follows:

$$f_E \circ g_E : [I_s, I_{\hat{d}}] \rightarrow [u, v, d]^T \in \mathbb{R}^3, \text{ where } g_E : I \in \mathbb{R}^{h \times w \times N} \rightarrow I \in \mathbb{R}^{H \times W \times N}. \quad (4)$$

124 Furthermore, we train the framework in a twin architecture fashion as proposed in [24, 25, 11, 20, 16,  
125 17, 21, 22] on the KeypointNet task.

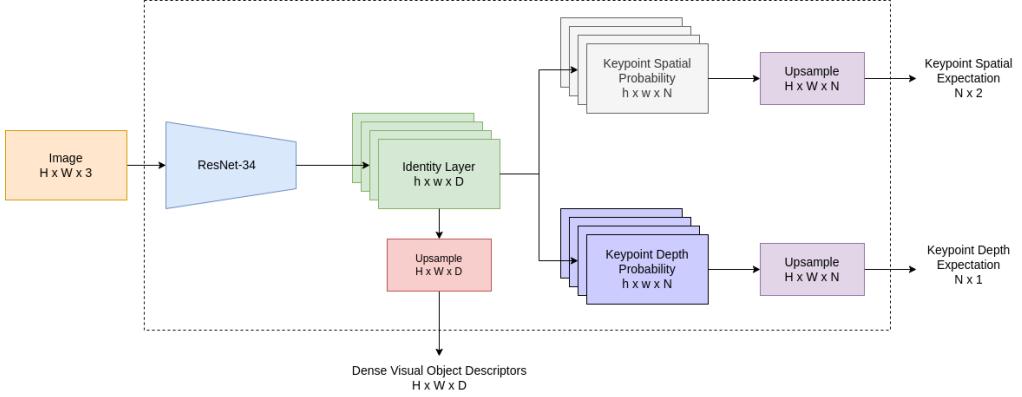


Figure 2: Illustration of novel framework designed to efficiently compute and seamlessly extract dense visual object descriptors. During inference we extract dense visual object descriptors directly from the network and ignore predicted spatial-depth expectation of the keypoints.

### 126 3.3 Loss Functions

127 For training, we directly adopt silhouette consistency loss ( $\mathcal{L}_{obj}$ ), variance loss ( $\mathcal{L}_{var}$ ) and separation  
128 loss ( $\mathcal{L}_{sep}$ ) functions from [28] to train the network on the keypoint prediction task. However, we  
129 modify the multi-view consistent loss and relative pose estimation loss. In the case of multi-view  
130 consistency loss we project the predicted spatial-depth expectation using camera intrinsics as follows:  
131

$$X_{cam} \in \mathbb{R}^{3 \times 1} = \mathcal{I}_{cam}^{-1} [u, v, 1.0]^T \otimes d, \text{ where } \mathcal{I}_{cam} \in \mathbb{R}^{3 \times 3} \text{ and } u, v, d \in \mathbb{R}^+. \quad (5)$$

132 Furthermore, we project the camera coordinates of the keypoints from one camera viewpoint to  
133 another camera viewpoint using relative transformation supplied from the synthetic augmentation  
134 procedure as follows:

$$\mathcal{L}_{mvc} \in \mathbb{R} = \mathcal{H}(\hat{X}_{cam}^B, \mathcal{T}_{A \rightarrow B} \hat{X}_{cam}^A), \text{ where } \hat{X}_{cam} = [X_{cam}, 1.0]^T \in \mathbb{R}^{4 \times 1}, \quad (6)$$

135 In Equation 6,  $\mathcal{T}_{A \rightarrow B} \in SE(3) \in \mathbb{R}^{4 \times 4}$  is a Special Euclidean Group [35] which is relative trans-  
136 formation from camera-frame A to camera-frame B. We use Huber loss  $\mathcal{H}$  as it produces smoother  
137 gradients for framework optimization. Furthermore, we do not discard the relative transformation  
138 information to calculate the realative pose loss as suggested in [28] and being influenced from [29]  
139 we modified the relative pose loss as follows:

$$\mathcal{L}_{pose} = \|\log(\mathcal{T}_{truth}^\dagger \mathcal{T}_{pred})\|, \text{ where } \log : SE(3) \rightarrow \mathfrak{se}(3) \text{ and } \mathcal{T}^\dagger = \begin{bmatrix} R^T & -R^T t \\ 0^T & 1 \end{bmatrix} \in SE(3). \quad (7)$$

### 140 3.4 Robot Grasping Pipeline

141 To use the proposed framework as a robot grasping pipeline, we extract dense visual object descriptors  
142 from the network and store one single descriptor of objects in a database manually for now. During

143 inference, we extract dense visual object descriptors from the network and query the descriptor from  
 144 the database to find the closest match as follows:

$$\mathbb{E}[u^*, v^*]_d = \operatorname{argmin}_{u, v} \exp - \left( \frac{\|I_D[u, v] - d\|}{\exp(t)} \right)^2 \quad (8)$$

145 Where  $t \in \mathbb{R}$  controls the kernel width influencing the search space to compute the optimal spatial  
 146 expectation  $\mathbb{E}[u^*, v^*]_d$  of the query descriptor  $d \in \mathbb{R}^D$  in the descriptor image  $I_D \in \mathbb{R}^{H \times W \times D}$ . The  
 147 computed spatial expectation is projected to robot frame using camera intrinsics and pose to perform  
 148 a pinch grasp. Furthermore, Franka Emika 7-DOF robot manipulator with 2 jaw gripper and wrist  
 149 mounted Intel Realsense D435 camera is used as testing setup as illustrated in Figure 3.

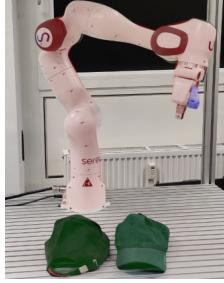


Figure 3: Illustration of the robot grasping pipeline setup. In the image, the robot is highlighted in red, the caps are highlighted in green and the camera is highlighted in blue.

## 150 4 Experiments & Results

### 151 4.1 Dense Object Nets

152 We implemented our training and benchmarking using “PyTorch-Lightning”[36] and “PyTorch”[37]  
 153 libraries. Furthermore, we employ ADAM[38] optimizer to optimize the model for 2500 epochs  
 154 with learning rate of  $\alpha = 3 \times 10^{-4}$ ,  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$  with weight decay  $\eta = 10^{-4}$  to  
 155 benchmark the DON with Pixelwise NT-Xent loss as in [17] with a fixed batch size of 1 and 128  
 156 image-pair correspondences. As per the benchmarking results in Table 2, the robustness of the  
 157 descriptor increases as the dimension of the descriptor gets longer.

Table 2: Benchmark of DON framework for GPU consumption and AUC for  $PCK@k$ ,  $\forall k \in [1, 100]$   
 metric.

DON benchmark				
Descriptor Size ( $D$ )	3	8	4	32
AUC for $PCK@k$	$0.922 \pm 0.006$	$0.933 \pm 0.011$	$0.948 \pm 0.012$	$0.953 \pm 0.008$
VRAM Usage (GB)	9.377	13.717	20.479	30.067

158 The AUC for  $PCK@k$ ,  $\forall k \in [1, 100]$  is computed with 256 image-pair correspondences and the  
 159 metric’s mean and std. deviation is calculated from benchmarking 3 DON models trained for a  
 160 single descriptor dimension. We could not train the descriptor dimension of 64 and 128 due to the  
 161 limited VRAM. Furthermore, to inspect the results of trained DON, a interface is built using the  
 162 PyGame library [39] to visualize the results of the trained DON. The mouse pointer in image space is  
 163 mapped to the pixel and the descriptor at that pixel is queried in another image-descriptor space.  
 164 We further use the spatial probability of the descriptor to visualize the queried descriptor in the image  
 165 space using Equation 8 to identify if there are any multi-modal spatial activations in the descriptor  
 166 spaces and there are none as shown in Figure 4.

### 167 4.2 Our Framework

168 To train our framework we employ ADAM optimizer to optimize the model for 2500 epochs with  
 169 learning rate of  $\alpha = 1 \times 10^{-3}$ ,  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$  with no weight decay. We further use a

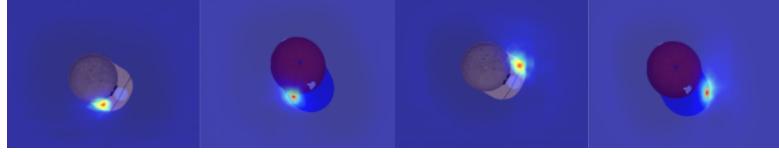


Figure 4: Depiction of the spatial probability heatmaps of the descriptor in the image space. We set the temperature in the Equation 8 to 1.1 and render the spatial probability heatmaps in the interface. The first and second image from the left and the right highlights the semantically equivalent descriptors in the image space.

170 fixed batch size of 1 and use StepLR scheduler with step size of 2500 and gamma of 0.9 to train the  
 171 model. At first, we trained our model with 16 keypoints with margin of 10 pixels as a hyperparameter  
 172 for the separation loss and later we trained the models with 128 keypoints with margin of 2 pixels.

Table 3: Benchmark of our framework for GPU consumption and AUC for  $PCK@k$ ,  $\forall k \in [1, 100]$  metric.

Our framework with 16 keypoints				
Descriptor Size ( $D$ )	64	128	256	512
AUC for $PCK@k$	$0.922 \pm 0.006$	$0.933 \pm 0.011$	$0.948 \pm 0.012$	$0.953 \pm 0.008$
VRAM Usage (GB)	3.799	4.191	5.241	7.341
Our framework with 128 keypoints				
Descriptor Size ( $D$ )	64	128	256	512
AUC for $PCK@k$	$0.922 \pm 0.006$	$0.933 \pm 0.011$	$0.948 \pm 0.012$	$0.953 \pm 0.008$
VRAM Usage (GB)	4.913	5.409	6.551	7.915
Our framework with 128 keypoints				
Descriptor Size ( $D$ )	3	16	32	64
AUC for $PCK@k$	$0.922 \pm 0.006$	$0.933 \pm 0.011$	$0.948 \pm 0.012$	$0.953 \pm 0.008$
VRAM Usage (GB)	4.913	5.409	6.551	7.915

### 173 4.3 Robot Grasping Pipeline

174 For the robot grasping pipeline, we trained our framework with actual caps. As the synthetic data  
 175 generation only needs mask and depth information, we could create mask in no time. Additionally,  
 176 while training the framework we do not need the actual real world depth information as the framework  
 177 computes it's own. We later extracted the dense visual local descriptors from the framework and  
 178 visually inspected for any inconsistencies in the descriptor space as shown in Figure 5. and found  
 179 it to be consistent. Furthermore, we did not use the models trained on the synthetic dataset as the  
 180 representations were inconsistent with the real caps.



Figure 5: Visual inspection of the dense visual descriptors space of the real caps.

181 For robot grasping, a descriptor is picked from the descriptor space and queried in the real-time such  
 182 that robot can pinch grasp the object. We could successfully grasp the caps with the robot as shown  
 183 in Figure 6.

184 As our framework inertly regresses keypoints on the object, we could use it as an alternative approach  
 185 to grasp the caps by computing the pose generated by the keypoints considering the actual depth  
 186 information instead of network regressed depth information. We extract the spatial probability of  
 187 each keypoint from the framework and deactivate spatial probabilities where the depth information is

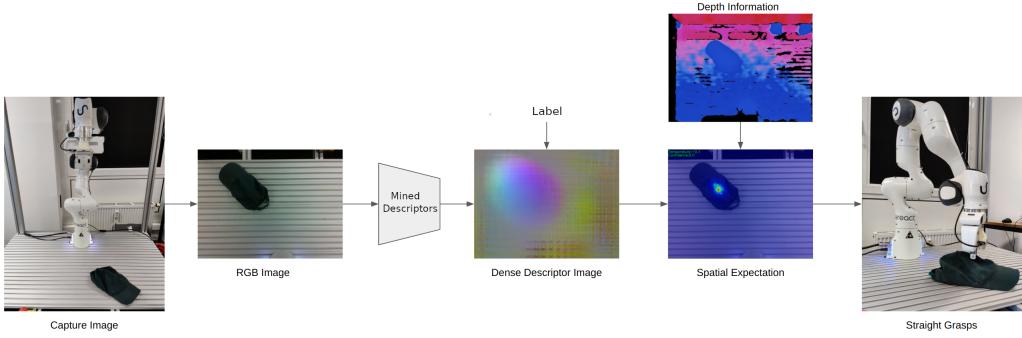


Figure 6: Depiction of the straight robot grasping pipeline.

188 missing as the depth image from the camera is noisy. Furthermore, the spatial expectations of the  
 189 keypoints is projected to camera frame to calculate a 6D pose in the camera frame. The 6D pose is  
 190 transformed the robot frame to perform a aligned grasp as shown in Figure 7.

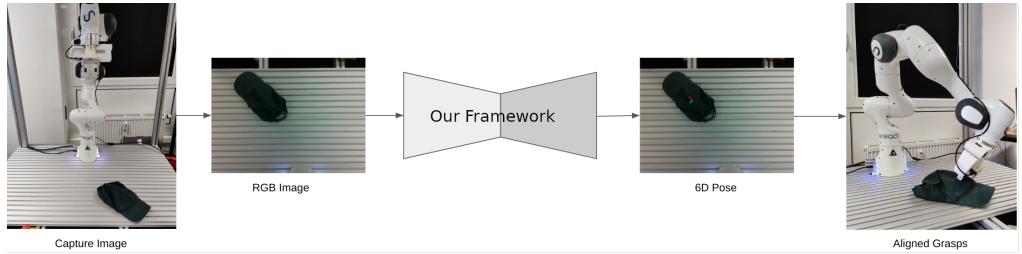


Figure 7: Illustration of the aligned robot grasping pipeline.

191 We did not evaluate the robot grasping pipeline.

## 192 5 Conclusion

193 We present a novel framework for mining dense visual object descriptors without explicitly training  
 194 DON. By leveraging synthetic augmentation data generation and a novel deep learning architecture,  
 195 our approach produces robust and denser visual local descriptors while consuming XX% lesser  
 196 computational resources than originally proposed framework. Futhermore, it eliminates the additional  
 197 computation of mapping large number of image-pair correspondence mapping. We demonstrate the  
 198 application of our framework as a robot-grasping pipeline in two methodolgies one of which our  
 199 framework demonstrates its capabilities to produce object specific 6D poses for robot grasping.

## 200 6 Future Work

201 The framework will be extended to produce multi-object dense visual descriptors in corporating  
 202 cluttered scenes. Furthermore, we will start incorporating ROI layers in the framework and adding  
 203 additional object classification loss function to the framework.

## 204 Broader Impact

## 205 References

- 206 [1] N. Blomkamp, H. Zimmer, and S. Kinberg. *Chappie*. Sony Pictures Home Entertainment, 2015.
- 207 [2] G. Lucas and S. Ulstein. *Star wars*. 20th Century-Fox, 1977.
- 208 [3] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran,  
 209 T. Graepel, et al. “A general reinforcement learning algorithm that masters chess, shogi, and Go through  
 210 self-play”. In: *Science* 362.6419 (2018), pp. 1140–1144.

- 211 [4] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I.  
 212 Antonoglou, V. Panneershelvam, M. Lanctot, et al. “Mastering the game of Go with deep neural networks  
 213 and tree search”. In: *nature* 529.7587 (2016), pp. 484–489.
- 214 [5] B. Entertainment. *World of warcraft*. Insight Editions, Division Of Palac, 2013.
- 215 [6] A. Krizhevsky, I. Sutskever, and G. E. Hinton. “Imagenet classification with deep convolutional neural  
 216 networks”. In: *Communications of the ACM* 60.6 (2017), pp. 84–90.
- 217 [7] J. Long, E. Shelhamer, and T. Darrell. “Fully convolutional networks for semantic segmentation”. In:  
 218 *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 3431–3440.
- 219 [8] K. He, G. Gkioxari, P. Dollár, and R. Girshick. “Mask r-cnn”. In: *Proceedings of the IEEE international  
 220 conference on computer vision*. 2017, pp. 2961–2969.
- 221 [9] R. A. Güler, N. Neverova, and I. Kokkinos. “Densepose: Dense human pose estimation in the wild”. In:  
 222 *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 7297–7306.
- 223 [10] sereact. “AI Robots for Production. Today”. <https://sereact.ai/en>.
- 224 [11] P. R. Florence, L. Manuelli, and R. Tedrake. “Dense object nets: Learning dense visual object descriptors  
 225 by and for robotic manipulation”. In: *arXiv preprint arXiv:1806.08756* (2018).
- 226 [12] P. Sundaresan, J. Grannen, B. Thananjeyan, A. Balakrishna, M. Laskey, K. Stone, J. E. Gonzalez, and K.  
 227 Goldberg. “Learning Rope Manipulation Policies Using Dense Object Descriptors Trained on Synthetic  
 228 Depth Data”. In: *CoRR* abs/2003.01835 (2020). arXiv: 2003.01835.
- 229 [13] C.-Y. Chai, K.-F. Hsu, and S.-L. Tsao. “Multi-step Pick-and-Place Tasks Using Object-centric Dense  
 230 Correspondences”. In: *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems  
 231 (IROS)*. 2019, pp. 4004–4011. DOI: 10.1109/IROS40897.2019.8968294.
- 232 [14] P. Florence, L. Manuelli, and R. Tedrake. “Self-supervised correspondence in visuomotor policy learning”.  
 233 In: *IEEE Robotics and Automation Letters* 5.2 (2019), pp. 492–499.
- 234 [15] A. Ganapathi et al. “Learning Dense Visual Correspondences in Simulation to Smooth and Fold Real  
 235 Fabrics”. In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*. 2021, pp. 11515–  
 236 11522. DOI: 10.1109/ICRA48506.2021.9561980.
- 237 [16] A. Kupcsik, M. Spies, A. Klein, M. Todescato, N. Wanek, P. Schillinger, and M. Bürger. “Supervised  
 238 Training of Dense Object Nets using Optimal Descriptors for Industrial Robotic Applications”. In: *arXiv  
 239 preprint arXiv:2102.08096* (2021).
- 240 [17] D. B. Adrian, A. G. Kupcsik, M. Spies, and H. Neumann. “Efficient and Robust Training of Dense  
 241 Object Nets for Multi-Object Robot Manipulation”. In: *2022 International Conference on Robotics and  
 242 Automation (ICRA)*. IEEE. 2022, pp. 1562–1568.
- 243 [18] R. Hartley and A. Zisserman. *Multiple view geometry in computer vision*. Cambridge university press,  
 244 2003.
- 245 [19] M. Belkin and P. Niyogi. “Laplacian eigenmaps for dimensionality reduction and data representation”. In:  
 246 *Neural computation* 15.6 (2003), pp. 1373–1396.
- 247 [20] P. R. Florence. “Dense visual learning for robot manipulation”. PhD thesis. Massachusetts Institute of  
 248 Technology, 2020.
- 249 [21] D. Hadjivelichkov and D. Kanoulas. “Fully Self-Supervised Class Awareness in Dense Object Descrip-  
 250 tors”. In: *5th Annual Conference on Robot Learning*. 2021.
- 251 [22] L. Yen-Chen, P. Florence, J. T. Barron, T.-Y. Lin, A. Rodriguez, and P. Isola. *NeRF-Supervision: Learning  
 252 Dense Object Descriptors from Neural Radiance Fields*. 2022. DOI: 10.48550/ARXIV.2203.01913.
- 253 [23] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. “Nerf: Representing  
 254 scenes as neural radiance fields for view synthesis”. In: *Communications of the ACM* 65.1 (2021), pp. 99–  
 255 106.
- 256 [24] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. “A simple framework for contrastive learning of visual  
 257 representations”. In: *International conference on machine learning*. PMLR. 2020, pp. 1597–1607.
- 258 [25] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. “Barlow twins: Self-supervised learning via  
 259 redundancy reduction”. In: *International Conference on Machine Learning*. PMLR. 2021, pp. 12310–  
 260 12320.
- 261 [26] C.-Y. Chai, K.-F. Hsu, and S.-L. Tsao. “Multi-step pick-and-place tasks using object-centric dense  
 262 correspondences”. In: *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*.  
 263 IEEE. 2019, pp. 4004–4011.
- 264 [27] M. E. Fathy, Q.-H. Tran, M. Z. Zia, P. Vernaza, and M. Chandraker. “Hierarchical metric learning and  
 265 matching for 2d and 3d geometric correspondences”. In: *Proceedings of the european conference on  
 266 computer vision (ECCV)*. 2018, pp. 803–819.
- 267 [28] S. Suwajanakorn, N. Snavely, J. J. Tompson, and M. Norouzi. “Discovery of latent 3d keypoints via  
 268 end-to-end geometric reasoning”. In: *Advances in neural information processing systems* 31 (2018).
- 269 [29] W. Zhao, S. Zhang, Z. Guan, W. Zhao, J. Peng, and J. Fan. “Learning deep network for detecting 3d  
 270 object keypoints and 6d poses”. In: *Proceedings of the IEEE/CVF Conference on computer vision and  
 271 pattern recognition*. 2020, pp. 14134–14142.

- 272 [30] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song,  
273 H. Su, et al. “Shapenet: An information-rich 3d model repository”. In: *arXiv preprint arXiv:1512.03012*  
274 (2015).
- 275 [31] M. Denninger, M. Sundermeyer, D. Winkelbauer, Y. Zidan, D. Olefir, M. Elbadrawy, A. Lodhi, and  
276 H. Katam. “Blenderproc”. In: *arXiv preprint arXiv:1911.01911* (2019).
- 277 [32] S. Marcel and Y. Rodriguez. “Torchvision the machine-vision package of torch”. In: *Proceedings of the*  
278 *18th ACM international conference on Multimedia*. 2010, pp. 1485–1488.
- 279 [33] W. Kabsch. “A solution for the best rotation to relate two sets of vectors”. In: *Acta Crystallographica*  
280 *Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography* 32.5 (1976), pp. 922–  
281 923.
- 282 [34] K. He, X. Zhang, S. Ren, and J. Sun. “Deep residual learning for image recognition”. In: (2016), pp. 770–  
283 778.
- 284 [35] W. P. Thurston. “Three-Dimensional Geometry and Topology, Volume 1”. In: *Three-Dimensional Geome-*  
285 *try and Topology, Volume 1*. Princeton university press, 2014.
- 286 [36] W. A. Falcon. “Pytorch lightning”. In: *GitHub 3* (2019).
- 287 [37] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein,  
288 L. Antiga, et al. “Pytorch: An imperative style, high-performance deep learning library”. In: *Advances in*  
289 *neural information processing systems* 32 (2019).
- 290 [38] D. P. Kingma and J. Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint*  
291 *arXiv:1412.6980* (2014).
- 292 [39] S. Kelly. “Basic introduction to pygame”. In: *Python, PyGame and Raspberry Pi Game Development*.  
293 Springer, 2016, pp. 59–65.