

Training Dense Object Nets: A Novel Approach

1st Kanishk Navale
Sereact GmbH.
Stuttgart, Germany
kanishk.navale@sereact.ai

2nd Ralf Gulde
Sereact GmbH.
Stuttgart, Germany
ralf.gulde@sereact.ai

3rd Marc Tuscher
Sereact GmbH.
Stuttgart, Germany
marc.tuscher@sereact.ai

4th Oliver Riedel
ISW - Universität Stuttgart
Stuttgart, Germany
oliver.riedel@isw.uni-stuttgart.de

Abstract—Our work proposes a novel framework that addresses the computational limitations associated with training Dense Object Nets (DON) while achieving robust and dense visual object descriptors. DON’s descriptors are known for their robustness to viewpoint and configuration changes, but training these requires image pairs with computationally expensive correspondence mapping. This limitation hampers dimensionality and robustness, thereby restricting object generalization. To overcome this, we introduce data generation procedure using synthetic augmentation and a novel deep learning architecture that produces denser visual descriptors with reduced computational demands. Notably, our framework eliminates the need for image pair correspondence mapping and showcases its application in a robotic grasping pipeline. Experimental results demonstrate that our approach yields descriptors as robust as those generated by DON.

Index Terms—Dense Object Nets, robot grasping, generalized object representation, reduced computation costs.

I. INTRODUCTION

Significant strides have been made in robotics, including self-driving cars and humanoid robots with advanced capabilities. Integration of AI models like ChatGPT [1] and PaLM [2] enhances human-robot interactions and object perception. However, deploying Large Language Models (LLMs) such as ChatGPT and Palm-E poses challenges. Their size and complexity require substantial computational resources, raising concerns about energy consumption, environmental impact, and the digital divide [3, 4]. Responsible resource allocation and addressing ethical implications are crucial for balancing progress and sustainability.

Currently, in robotics, typical industrial robots perform repetitive operations based on pre-programmed instructions, finding the ideal object representation for grasping and manipulation tasks still needs to be answered. Existing representations may be unable to understand an object’s geometrical and structural information, rendering them unsuitable for complex tasks. In recent work, Florence et al. [5] introduced a novel visual object representation termed “dense visual object descriptors” to the robotics community. This representation, generated by the Dense Object Nets (DON) framework, converts each pixel in an image ($I[u, v] \in \mathbb{R}^3$) into a higher-dimensional embedding ($I_D[u, v] \in \mathbb{R}^D$) such that $D \in \mathbb{N}^+$, using image-pair correspondences as input. These dense visual object descriptors provide a generalized representation of objects to a certain extent.

The DON framework has shown promise in various domains, including rope manipulation [6], block manipulation [7], robot

control [8], fabric manipulation [9], and robot grasp pose estimation [10, 11]. Adrian et al. [11] demonstrated that DON can be trained on synthetic data and still generalize well to real-world objects. Furthermore, Adrian et al. [11] highlight the importance of the dimensionality of the embedding in determining the quality of the descriptors produced by the DON framework.

In this paper, we address the challenge posed by the computationally intensive nature of DON and propose a new framework for training DON in a computationally efficient manner. Furthermore, we introduce a novel synthetic data generation pipeline that generates a complete dataset from one image and mask pair. Additionally, the synthetic data generation pipeline does not rely on the noisy depth information produced by today’s consumer-grade depth cameras. We also demonstrate one of the applications of our framework as a robotic grasping pipeline. Our approach aims to contribute to developing a sustainable and efficient, and economical solution for industrial robotics applications.

II. RELATED WORKS

Florence et al. [5] introduced the Pixelwise Contrastive loss function to train DON, which involves sampling pixels in an image-pair and computing the Contrastive loss between the pixels in the first image and those in the second image. This optimization procedure aims to improve the descriptor based on a similarity metric. However, the Pixelwise Contrastive loss function is computationally expensive and requires numerous matching and non-matching image-pair correspondences to work optimally. When optimizing DON using a large number (N) of image-pair correspondences, the computational resources consumed by the optimization procedure increase significantly due to the exponential growth of pixelwise descriptor similarity comparisons (2^N).

In their work, Florence [12] discovered that the Pixelwise Contrastive loss function used to train DON might yield poor performance if a computed correspondence is spatially inconsistent. They also highlighted that the precision of contrastive-trained models could be sensitive to the relative weighting between positive and negative sampled pixels. To address these limitations, Florence proposed a new continuous sampling-based loss function called the Pixelwise Distribution loss. This novel loss function leverages smooth and continuous pixel space sampling instead of the discrete pixel space sampling method employed by the Pixelwise Contrastive loss.

The Pixelwise Distribution loss eliminates the need for non-matching correspondences, leading to significant savings in computation resources.

On a different note, Kupcsik et al. [10] utilized Laplacian Eigenmaps [13] to embed a 3D object model into an optimally generated embedding space, serving as the target for training DON in a supervised fashion. However, this methodology does not reduce the computational resource consumption required to train DON. In contrast, Hadjivelichkov and Kanoulas [14] employed offline unsupervised clustering based on confidence in object similarities to generate hard and soft correspondence labels. These labels were then used as matching and non-matching correspondences to train DON effectively.

Building upon the concept of SIMCLR-inspired frameworks [15, 16], Adrian et al. [11] introduced a similar architecture and another novel loss function called the Pixelwise NTXent Loss. This loss function robustly trains DON by leveraging synthetic correspondences computed from image augmentations and non-matching image correspondences. Notably, Adrian et al.’s experiments demonstrated that the novel loss function is invariant to batch size variations, unlike the Pixelwise Contrastive Loss. Furthermore, it is worth noting that most of the discussed optimization methodologies heavily rely on correspondences to train DON effectively.

Moving on to the aspect of image-pair correspondences and dataset engineering, the DON training strategy proposed in [5, 12] relies on depth information to compute correspondences between image pairs using camera intrinsics and pose information [17]. However, when utilizing consumer-grade depth cameras to capture depth information, the resulting depth data can be noisy, particularly when dealing with tiny, reflecting objects common in industrial environments. Noisy depth information hampers the computation of consistent spatial correspondences in an image pair. To overcome this challenge, Kupcsik et al. [10] associated 3D models of objects with image views, effectively training DON without relying on depth information. Their approach proved efficient for smaller, texture-less, and reflective objects. Additionally, Kupcsik et al. compared different training strategies for producing 6D grasps on industrial objects and demonstrated that a unique supervised training approach enhances pick-and-place resilience in industry-relevant tasks.

In contrast, Yen-Chen et al.[18] employed NeRF[19], a method that reconstructs a 3D scene from a sequence of images captured by a smartphone camera. They extracted correspondences from the synthetically reconstructed scene to train DON. Remarkably, Adrian et al.’s experiments indicated that DON trained on synthetic data generalizes well to real-world objects. Furthermore, they adopted the $PCK@k$ metric, commonly used in [20, 21], to evaluate and benchmark DON’s performance in cluttered scenes that were previously not extensively studied.

On further exploration of frameworks that could generalize objects, we ended up at the framework introduced by Suwanjanakorn et al. [22]. The authors presented a framework to predict geometrically consistent keypoints. These keypoints possess the capability to generalize objects. Upon further

investigation, we discovered that one of the layers within the framework bears a resemblance to dense visual object descriptors. This similarity is attributed to the inherent property of the framework, which involves regressing keypoints that hold semantic equivalence across objects. Building upon the framework introduced in [22], Zhao et al. [23] extended it to a multi-object scene.

In our work, we do not use any loss functions as proposed in [5, 12, 10, 11, 14, 18] to train DON. However, we adopt the network architecture from DON [5] as our architecture’s backbone and train on the task of the KeypointNet[22, 23] with few network modifications. Moreover, we evaluate the descriptor’s robustness produced by our framework on the $PCK@k$ metric as in comparision to benchmarks in [11] as it is the only benchmark available for DON. Furthermore, we compare the computational resource consumption used for training both frameworks.

III. METHODOLOGY

In this section, we outline the methodologies employed. Our approach encompasses synthetic dataset engineering, a novel framework, loss function modifications and a comprehensive grasping pipeline.

Firstly, we focus on synthetic dataset engineering accommodating spatial, colour and background augmentation. The colour and background augmentations help the framework to predict object-oriented descriptors. Secondly, we present a novel framework designed to reduce computational resource consumption with loss function modifications to optimize performance. Lastly, we introduce a comprehensive robot grasping pipeline exploiting the generalizing capabilities of our framework.

A. Dataset Engineering

We have chosen the cap object for creating a synthetic dataset as the cap mesh models are readily available in the Shapenet library [24] as it contains rich object information, including textures. Furthermore, we choose five cap models from the Shapenet library and use Blenderproc [25] to generate the synthetic dataset. We save one RGB image, mask, and depth for each cap model from the synthetic scene. Additionally, we employ synthetic augmentations as proposed in [11] to synthetically spatial augment the cap’s position and rotation in an image, including background randomization using Torchvision [26] library. An augmented image pair is sampled randomly to generate camera poses for different viewpoints. Additionally, image-pair correspondences are computed ¹ as illustrated in the Figure 1. We only compute 24 image-pair correspondences for an image-pair as we found that 24 image correspondences yield stable computation for translation and rotation for all the objects. Using depth information, we project the computed correspondences to the camera frame and compute the relative transformation between two camera-frame coordinates of the correspondences using Kabsch’s transformation [27]. Moreover, mask and depth images are not used during inference.

¹GitHub Link:<https://shorturl.at/dimpF>



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.

B. Framework & Mining Strategy

As a backbone, we employ ResNet-34 architecture [28]. 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 identical 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 are 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 [22], 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:

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

and depth as follows:

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

We incorporate continuous sampling method f_E from [12, 22] to convert the upsampled predicted spatial-probability and depth of a keypoint to spatial-depth expectation as follows:

$$f_E \circ g_E : [I_s, I_{\hat{d}}] \rightarrow [u, v, d]^T \in \mathbb{R}^3, \quad (4)$$

where, $g_E : I \in \mathbb{R}^{h \times w \times N} \rightarrow I \in \mathbb{R}^{H \times W \times N}$. Furthermore, we train the framework in a twin architecture fashion as proposed in [15, 16, 5, 12, 10, 11, 14, 18] on the modified KeypointNet task.

C. Loss Functions

For training, we directly adopt silhouette consistency loss (\mathcal{L}_{obj}), variance loss (\mathcal{L}_{var}) and separation loss (\mathcal{L}_{sep}) functions from [22] to train the network on the keypoint prediction task. However, we modify the multi-view consistent loss and relative pose estimation loss. In the case of multi-view consistency loss we project the predicted spatial-depth expectation using camera intrinsics as follows:

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

where, $\mathcal{I}_{cam} \in \mathbb{R}^{3 \times 3}$ and $u, v, d \in \mathbb{R}^+$. Furthermore, we project the camera coordinates of the keypoints from one

camera viewpoint to another camera viewpoint using relative transformation supplied from the synthetic augmentation procedure as follows:

$$\mathcal{L}_{mvc} \in \mathbb{R} = \mathcal{H}(\hat{X}_{cam}^B, \mathcal{T}_{A \rightarrow B} \hat{X}_{cam}^A), \quad (6)$$

where, $\hat{X}_{cam} = [X_{cam}, 1.0]^T \in \mathbb{R}^{4 \times 1}$. In Equation 6, $\mathcal{T}_{A \rightarrow B} \in SE(3)$ is a Special Euclidean Group [29] which is relative transformation from camera-frame A to camera-frame B . We use Huber loss \mathcal{H} as it produces smoother gradients for framework optimization. Furthermore, we do not discard the relative transformation information to calculate the relative pose loss as suggested in [22]. Moreover, being influenced from [23] we modified the relative pose loss as follows:

$$\mathcal{L}_{pose} = \|\log(\mathcal{T}_{truth}^\dagger \mathcal{T}_{pred})\|, \quad (7)$$

where, $\log : SE(3) \rightarrow \mathfrak{se}(3)$ and $\mathcal{T}^\dagger = \begin{bmatrix} R^T & -R^T t \\ 0^T & 1 \end{bmatrix}$.

D. Robot Grasping Pipeline

To use the proposed framework as a robot grasping pipeline, we extract dense visual object descriptors from the network and store one single descriptor of objects in a database manually for now. During inference, we extract dense visual object descriptors from the network and query the descriptor from the database to find the closest match as follows:

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

where, $\|I_D[u, v] - d\| \in \mathbb{R}^{H \times W}$ Furthermore, $t \in \mathbb{R}$ controls the kernel width influencing the search space to compute the optimal spatial 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 computed spatial expectation is projected to the robot frame using camera intrinsics and poses to perform a pinch grasp. Furthermore, the Franka Emika 7-DOF robot manipulator with two jaw gripper and wrist-mounted Intel Realsense D435 camera is used as a testing setup as illustrated in Figure 3.

IV. EXPERIMENTS & RESULTS

In this section, we outline the benchmarking results employed from the methodologies. We benchmark the DON framework with Pixelwise NT-Xent loss as in [11] and our framework with our revision of the loss function on a 48GB VRAM GPU. We benchmark the descriptor's robustness with the $AUC \pm \sigma$ for $PCK@k, \forall k \in [1, 100]$ metric. Furthermore, we benchmark the computational resource consumption of the DON and our frameworks. We also demonstrate the application of our framework as a robot-grasping pipeline in two methodologies, one of which our framework demonstrates its capabilities to produce object-specific 6D poses for robot grasping.

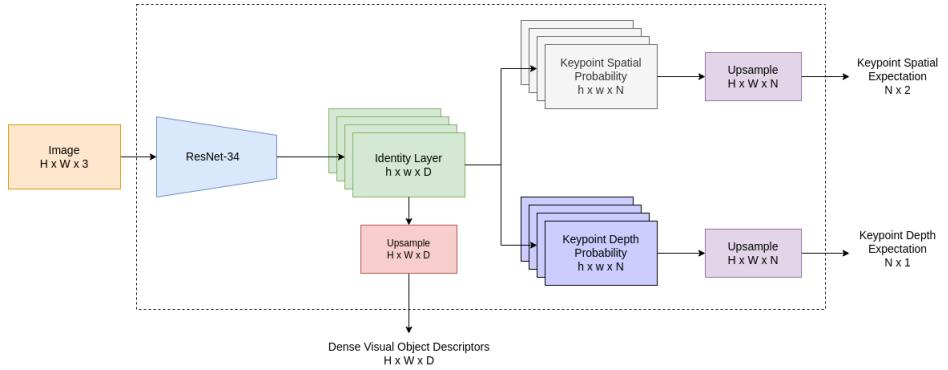


Figure 2. Illustration of the novel framework designed to compute and seamlessly extract dense visual object descriptors efficiently. During inference, we extract dense visual object descriptors directly from the network and ignore predicted spatial-depth expectations of the keypoints.



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

A. Training Setup

We implemented training and benchmarking using “PyTorch-Lightning”[30] and “PyTorch”[31] libraries. Furthermore, we employ ADAM[32] optimizer to optimize the model for 2500 epochs with a learning rate of $\alpha = 3 \times 10^{-4}$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$ with weight decay $\eta = 10^{-4}$ to benchmark the DON with Pixelwise NT-Xent loss as in [11] with a fixed batch size of 1 and 128 image-pair correspondences.

To train our framework, we employ an ADAM optimizer to optimize the model for 2500 epochs with a 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 fixed batch size of 1 and the StepLR scheduler with a step size 2500 and a gamma of 0.9 to train the model with all the loss weights to 1.0 except variance loss weight to 1×10^{-3} . We trained the three models with 128 keypoints with a margin of 2 pixels for each descriptor dimension. We specifically chose 128 keypoints as it aligns with the notion that DON is benchmarked with 128 image-pair correspondences.

B. Benchmarking & Results

The $AUC \pm \sigma$ for $PCK@k, \forall k \in [1, 100]$ is computed with 256 image-pair correspondences for both models. The metrics mean and std. deviation is calculated from benchmarking three models trained for each descriptor dimension. Due to the limited GPU VRAM capacity, we could not train DON for descriptor dimensions greater than 32. As per Table I, both frameworks benefit while training them for longer descriptor dimensions.

Furthermore, the higher values infer robust descriptors in Table I. We notice that our framework works robustly as the descriptor dimension gets longer as the metric difference between DON and our framework reduces.

Table I
BENCHMARKING OUTCOMES FOR DESCRIPTORS’ EVALUATION METRIC.

Benchmarking for $AUC \pm \sigma$ for $PCK@k, \forall k \in [1, 100]$		
Descriptor Size (D)	Dense Object Nets	Our framework
3	0.922 ± 0.006	0.914 ± 0.009
8	0.933 ± 0.011	0.928 ± 0.015
16	0.948 ± 0.012	0.945 ± 0.010
32	0.953 ± 0.008	0.950 ± 0.009
64	~	0.953 ± 0.006
128	~	0.957 ± 0.012
256	~	0.959 ± 0.008
512	~	0.962 ± 0.011

While training both frameworks, we monitor the GPU VRAM consumption. As per the benchmark results in Table II, the DON consumption increases as the descriptor dimensions get longer while our framework consumes a fraction of the computation resource. Furthermore, lower readings are better in Table II.

Table II
BENCHMARKING OUTCOMES FOR TRAINING COMPUTATION RESOURCE CONSUMPTION.

Benchmarking for GPU VRAM(GB) consumption		
Descriptor Size (D)	Dense Object Nets	Our framework
3	9.377	4.763
8	13.717	4.785
16	20.479	4.832
32	30.067	4.872
64	~	4.913
128	~	5.409
256	~	6.551
512	~	7.915

To check the impact of descriptors’ robustness compared to the number of keypoints, we trained our framework with 16 keypoints. Furthermore, we trained three additional models for each descriptor dimension 64, 128, 256, and 512. As per

the Table III compared to the results in Table I in page I, the descriptor's robustness decreased when the framework predicted 16 keypoints. Moreover, this reflects that number of keypoints in our framework and the number of image-pair correspondences in DON are directly proportional to the robustness of the descriptors.

Table III
BENCHMARK OF OUR FRAMEWORK WITH 16 KEYPOINTS FOR GPU VRAM(GB) CONSUMPTION AND $AUC \pm \sigma$ FOR $PCK@k, \forall k \in [1, 100]$ METRIC.

Our framework with 16 keypoints		
Descriptor Size (D)	Metric	VRAM Usage (GB)
64	0.948 ± 0.009	3.799
128	0.952 ± 0.010	4.191
256	0.955 ± 0.013	5.241
512	0.957 ± 0.006	7.341

C. Descriptor Inspection

Furthermore, to inspect the results of trained DON, an interface is built using the PyGame library [33] to visualize the results of the trained DON. The mouse pointer in the image space is mapped to the pixel, and the descriptor at that pixel is queried in another image-descriptor space. We further use the spatial probability of the descriptor to visualize the queried descriptor in the image space using Equation 8 in page 3. We identify if there are any multi-modal spatial activations in the descriptor spaces and none, as shown in Figure 4.



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.

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

D. Robot Grasping Pipeline

For robot grasping, a descriptor is picked from the descriptor space and queried in real-time such that robot can pinch-grasp the object. We could successfully grasp the caps with the robot, as shown in Figure 6. Furthermore, we did not evaluate the robot grasping for position and semantic object location offsets.

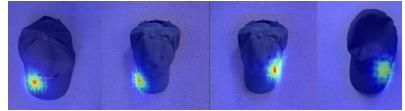


Figure 5. The image depicts the visual inspection of the dense visual descriptors space of the real caps using our developed interface. We trained our framework on the first two caps from the left, and our framework could generalize the object representations on an unseen cap while training illustrated in the first image from the right.

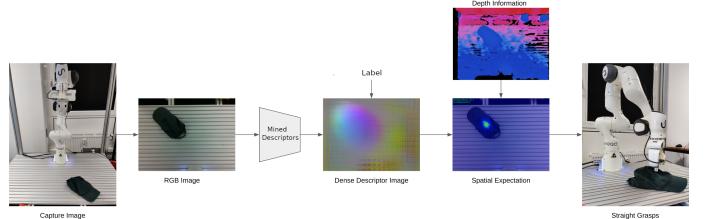


Figure 6. Depiction of the straight robot grasping pipeline.

As our framework inertly regresses keypoints on the object, we could use it as an alternative approach to grasp the caps by computing the pose generated by the keypoints considering the actual depth information instead of network-regressed depth information. We extract the spatial probability of each keypoint from the framework and deactivate spatial probabilities where the depth information is missing, as the depth image from the camera is noisy. Furthermore, the spatial expectations of the keypoints are projected to the camera frame to calculate a 6D pose in the camera frame. The 6D pose is transformed in the robot frame to perform an aligned grasp, as shown in Figure 7.

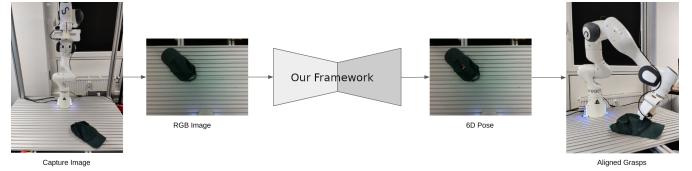


Figure 7. Illustration of the aligned robot grasping pipeline.

V. CONCLUSION

This paper introduces a novel framework for mining dense visual object descriptors without explicitly training DON. We have successfully eliminated the requirement for image-pair correspondence mapping in training DON by employing synthetic augmentation data generation and a novel deep-learning architecture. Our benchmarking results showcase the effectiveness of our framework in generating robust and denser visual local descriptors. However, it needs to outperform the original DON framework in robustness. Moreover, a notable advantage of our proposed framework is its significantly reduced computational resource consumption, amounting to a remarkable 86.67% decrease compared to the originally proposed framework. It is important to note that our current framework is limited to single object-dense visual descriptors.

Nevertheless, we have plans to extend our methodology to encompass the production of multi-object dense visual descriptors in cluttered scenes. By doing so, we aim to enhance the versatility and applicability of our framework in real-world scenarios. To demonstrate the practicality of our framework, we have integrated it into a robot-grasping pipeline using two distinct methodologies. Remarkably, our framework can generate object-specific 6D poses, enhancing robot grasping performance. This successful application further highlights the potential utility of our framework in real-world robotic systems.

REFERENCES

- [1] OpenAI. *GPT-4 Technical Report*. 2023. arXiv: 2303.08774 [cs.CL].
- [2] A. Chowdhery et al. “PaLM: Scaling Language Modeling with Pathways”. In: *arxiv:2204.02311* (2022).
- [3] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell. “On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?” In: *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 2021, pp. 610–623.
- [4] E. Strubell, A. Ganesh, and A. McCallum. “Energy and policy considerations for deep learning in NLP”. In: *arXiv preprint arXiv:1906.02243* (2019).
- [5] P. R. Florence, L. Manuelli, and R. Tedrake. “Dense object nets: Learning dense visual object descriptors by and for robotic manipulation”. In: *arXiv preprint arXiv:1806.08756* (2018).
- [6] P. Sundaresan, J. Grannen, B. Thananjeyan, A. Balakrishna, M. Laskey, K. Stone, J. E. Gonzalez, and K. Goldberg. “Learning Rope Manipulation Policies Using Dense Object Descriptors Trained on Synthetic Depth Data”. In: *CoRR* abs/2003.01835 (2020). arXiv: 2003.01835.
- [7] C.-Y. Chai, K.-F. Hsu, and S.-L. Tsao. “Multi-step Pick-and-Place Tasks Using Object-centric Dense Correspondences”. In: *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2019, pp. 4004–4011. DOI: 10.1109/IROS40897.2019.8968294.
- [8] P. Florence, L. Manuelli, and R. Tedrake. “Self-supervised correspondence in visuomotor policy learning”. In: *IEEE Robotics and Automation Letters* 5.2 (2019), pp. 492–499.
- [9] A. Ganapathi et al. “Learning Dense Visual Correspondences in Simulation to Smooth and Fold Real Fabrics”. In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*. 2021, pp. 11515–11522. DOI: 10.1109/ICRA48506.2021.9561980.
- [10] A. Kupcsik, M. Spies, A. Klein, M. Todescato, N. Waniek, P. Schillinger, and M. Bürger. “Supervised Training of Dense Object Nets using Optimal Descriptors for Industrial Robotic Applications”. In: *arXiv preprint arXiv:2102.08096* (2021).
- [11] D. B. Adrian, A. G. Kupcsik, M. Spies, and H. Neumann. “Efficient and Robust Training of Dense Object Nets for Multi-Object Robot Manipulation”. In: *2022 International Conference on Robotics and Automation (ICRA)*. IEEE. 2022, pp. 1562–1568.
- [12] P. R. Florence. “Dense visual learning for robot manipulation”. PhD thesis. Massachusetts Institute of Technology, 2020.
- [13] M. Belkin and P. Niyogi. “Laplacian eigenmaps for dimensionality reduction and data representation”. In: *Neural computation* 15.6 (2003), pp. 1373–1396.
- [14] D. Hadjivelichkov and D. Kanoulas. “Fully Self-Supervised Class Awareness in Dense Object Descriptors”. In: *5th Annual Conference on Robot Learning*. 2021.
- [15] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. “A simple framework for contrastive learning of visual representations”. In: *International conference on machine learning*. PMLR. 2020, pp. 1597–1607.
- [16] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. “Barlow twins: Self-supervised learning via redundancy reduction”. In: *International Conference on Machine Learning*. PMLR. 2021, pp. 12310–12320.
- [17] R. Hartley and A. Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [18] L. Yen-Chen, P. Florence, J. T. Barron, T.-Y. Lin, A. Rodriguez, and P. Isola. *NeRF-Supervision: Learning Dense Object Descriptors from Neural Radiance Fields*. 2022. DOI: 10.48550/ARXIV.2203.01913.
- [19] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. “Nerf: Representing scenes as neural radiance fields for view synthesis”. In: *Communications of the ACM* 65.1 (2021), pp. 99–106.
- [20] C.-Y. Chai, K.-F. Hsu, and S.-L. Tsao. “Multi-step pick-and-place tasks using object-centric dense correspondences”. In: *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2019, pp. 4004–4011.
- [21] M. E. Fathy, Q.-H. Tran, M. Z. Zia, P. Vernaza, and M. Chandraker. “Hierarchical metric learning and matching for 2d and 3d geometric correspondences”. In: *Proceedings of the european conference on computer vision (ECCV)*. 2018, pp. 803–819.
- [22] S. Suwajanakorn, N. Snavely, J. J. Tompson, and M. Norouzi. “Discovery of latent 3d keypoints via end-to-end geometric reasoning”. In: *Advances in neural information processing systems* 31 (2018).
- [23] W. Zhao, S. Zhang, Z. Guan, W. Zhao, J. Peng, and J. Fan. “Learning deep network for detecting 3d object keypoints and 6d poses”. In: *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*. 2020, pp. 14134–14142.
- [24] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al. “Shapenet: An information-rich 3d model repository”. In: *arXiv preprint arXiv:1512.03012* (2015).
- [25] M. Denninger, M. Sundermeyer, D. Winkelbauer, Y. Zidan, D. Olefir, M. Elbadrawy, A. Lodhi, and H. Katam. “Blenderproc”. In: *arXiv preprint arXiv:1911.01911* (2019).
- [26] S. Marcel and Y. Rodriguez. “Torchvision the machine-vision package of torch”. In: *Proceedings of the 18th ACM international conference on Multimedia*. 2010, pp. 1485–1488.
- [27] W. Kabsch. “A solution for the best rotation to relate two sets of vectors”. In: *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography* 32.5 (1976), pp. 922–923.
- [28] K. He, X. Zhang, S. Ren, and J. Sun. “Deep residual learning for image recognition”. In: (2016), pp. 770–778.
- [29] W. P. Thurston. “Three-Dimensional Geometry and Topology, Volume 1”. In: *Three-Dimensional Geometry and Topology, Volume 1*. Princeton university press, 2014.
- [30] W. A. Falcon. “Pytorch lightning”. In: *GitHub* 3 (2019).
- [31] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. “Pytorch: An imperative style, high-performance deep learning library”. In: *Advances in neural information processing systems* 32 (2019).
- [32] D. P. Kingma and J. Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014).
- [33] S. Kelly. “Basic introduction to pygame”. In: *Python, PyGame and Raspberry Pi Game Development*. Springer, 2016, pp. 59–65.