

Article

Not peer-reviewed version

Visual Imitation Learning from One-Shot Demonstration for Multi-Step Robot Pick-and-Place Tasks

[Shuang Lu](#)*, [Christian Haerdtlein](#), Johannes Schilp

Posted Date: 16 September 2024

doi: 10.20944/preprints202408.1123.v2

Keywords: Visual imitation learning; Robot; One-shot demonstration



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Article

Visual Imitation Learning from One-Shot Demonstration for Multi-Step Robot Pick-and-Place Tasks

Shuang Lu ^{1,*} , Christian Härdtlein ¹ and Johanness Schilp ^{1,2}

¹ Fraunhofer IGCV, Am Technologiezentrum 10, Augsburg, Germany

² Chair of Digital Manufacturing, Augsburg University, Am Technologiezentrum 8, Augsburg, Germany

* Correspondence: shuang.lu@igcv.fraunhofer.de

Abstract: Imitation learning, also known as programming by demonstration, has emerged as a method for intuitive robot programming by non-experts. Initially, the research field focused on physical interaction, where humans directly guided robots by manipulating their arms and tools to perform specific tasks. Over time, the field expanded to include observational learning, where robots learn from visual data. Visual imitation learning enables non-experts to demonstrate multi-step tasks by recording them as a single, continuous video. Existing approaches, particularly in household applications, often require large datasets to develop deep learning models that map visual inputs to robot actions. In contrast, this paper explores the use of one-shot demonstrations, which reduces data requirements and simplifies the programming process. This one-shot approach is particularly valuable for manufacturing scenarios, where reducing data collection time and training effort is essential. To facilitate this, the paper proposes a framework that maps hand trajectories to the robot's end effector. It consists of four essential components: hand detection, object detection, segmentation of trajectories into elementary skills, and skill learning. Each component is implemented and evaluated using recorded videos to demonstrate the effectiveness of the framework.

Keywords: visual imitation learning; robot; one-shot demonstration

1. Introduction

Recent changes in market trends, including increased emphasis on product differentiation, personalization, and high-mix, low-volume production, have highlighted the importance of flexibility in the manufacturing industry [1]. The wide variety of products being manufactured requires flexible robot systems, particularly for pick & place tasks. They must handle items of varying shapes, sizes, and weights. This adaptability allows manufacturers to respond quickly to changing production demands. However, the traditional robot programming approach is not feasible for flexible pick & place systems, as it requires expert knowledge and is time-consuming [2,3].

Learning from demonstration is a framework that allows robots to learn to perform tasks from human demonstrations of those tasks. This approach is intuitive for non-expert users, making robot programming more accessible. The demonstration can be performed by physically guiding a robot manipulator [4] or by recording it with a camera [5]. Recently, visual imitation learning gains attention thanks to the advancements in computer vision. This paradigm enables the learned task model to be independent of the robot manipulator's physics constraints while integrating of contextual information from the task environment.

Existing visual imitation learning approaches primarily focus on applications in household scenarios. The main characteristics of household environments include the availability of training data and the availability of open source virtual environments for data generation [6–9]. These factors allow for the training of large models that map visual inputs to robot actions. However, applications in manufacturing scenarios are more diverse and domain-specific. Unlike household tasks, where datasets such as those for kitchen applications are readily available [10], manufacturing lacks large,

standardized datasets. In addition, creating various virtual environments for data collection is not efficient in these cases, as the tasks and objects involved are often highly specialized. As a result, more efficient methods for training and demonstration are necessary to address these challenges in industrial contexts. This work aims to implement one-shot demonstration as a solution, focusing on efficient data collection. Human demonstrations are recorded using an RGB-D camera, which is widely employed in robotic applications [5,11,12] for capturing both color and depth information. To implement one-shot imitation learning for multi-step pick & place tasks, a framework consisting of four components is proposed: hand detection, object detection, segmentation of trajectories into elementary skills, and skill learning. The component-based framework makes it possible to integrate pre-trained models for tasks like hand and object detection.

The main contribution of this work is in the integration of the methods for each component into a coherent framework that enables effective one-shot visual imitation learning. For hand detection, two widely applied methods [13,14] are discussed and evaluated. In manufacturing scenarios, CAD files are often available. This approach reduces the need for extensive manual data collection by leveraging CAD models. Specifically, the approach by Denninger *et al.* [15] is used as an automatic pipeline to generate training data. For trajectory segmentation, an algorithm is proposed to divide multi-step hand movements into distinct actions. Dynamic Movement Primitives (DMPs) [16] are utilized for skill learning, as they offer adaptability to varying start and goal positions [11].

2. Related Work

To understand and review what is learned through visual imitation learning, the outcomes of visual imitation learning are categorized into three levels:

- **Skill Level:** This involves learning specific movements or actions that a robot must perform. Skills are the fundamental building blocks, such as grasping an object.
- **Task Level:** At this level, the robot learns to combine multiple skills to perform a more complex activity.
- **Goal Level:** This highest level involves understanding the overall objective of tasks.

At the **skill level**, the focus is on the acquisition of individual motor skills or actions for basic manipulations. Finn *et al.* presented a visual imitation learning method that enables a robot to learn new skills such as *push* and *place* from raw pixel input [6]. The policy is represented by Convolutional Neural Networks (CNN). The policy observation includes both the RGB image and the robot's joint angles and end-effector pose. A policy π is learned to map observations to robot actions. The approach integrates meta-learning with imitation learning, allowing a robot to reuse experience and quickly learn new skills from a single demonstration. This is achieved through a two-level learning process: the inner loop makes task-specific adjustments, while the outer loop updates meta-parameters across tasks. However, the two-level training process is complex and the learned model is adapting poorly to environment changes. Xin *et al.* developed a so-called IRMT-Net to predict the interaction region and motion trajectory from RGB-D images [9]. The input of the developed model consists of motion trajectories, motion category and cropped object images. To generate motion trajectory from RGB-D images, the hand detection method proposed by Shan *et al.* [13] is applied to extract bounding box from each RGB frame. The center coordinate of the bounding box and its depth value are taken as 3D hand coordinates. Faster-RCNN [17] is used to detect objects. The cropped image for the detected object is obtained based on the bounding box with the highest detection score. To evaluate the proposed approach, a dataset is created by labeling the motion trajectories for videos in the Epic-kitchens dataset [7] with 9236 videos. In the end, skills such as *pull open drawer* and *take cup* in kitchen scenarios are learned. The annotated dataset contributes to the further development of visual imitation learning for kitchen scenarios. Wen *et al.* introduced a method for learning task trajectories using a single demonstration video, captured by a statically mounted Photoneo 3D camera [18]. This camera records at 10 Hz. It provides gray-scale and depth images, allowing for detailed observation of the working

environment. The method involves tracking target objects to generate trajectories, which demonstrated effective performance. A key factor contributing to the success of this approach is the use of relatively large objects, which minimizes occlusion by the human hand during the demonstration.

At the **task level**, the focus is on understanding sequences of actions. The raw data are typically hand trajectories captured by cameras or other types of motion sensors. Qiu et al. presented a system with observing human demonstrations by an ASUS RGB-D camera [5]. During demonstration, a human worker performs an object handling task wearing a colored hand glove. This hand glove improves the accuracy and robustness for hand detection. The hand pose is estimated based on a deep learning model trained by 3D input data [19]. The human demonstration is segmented by Hidden Markov Models (HMMs) into motion primitives called skills. They include *pick up*, *place* and *locate*. The skills are then represented by Dynamic Movement Primitives (DMPs) [20], which allows the generalization to new goal positions. Notably, there are no general consensus for defining the semantic of skills in the existing works. Qiu et al. consider *pick up*, *place* and *locate* as skills [5]. However, Kyrarini et al. define them as *start arm moving*, *object grasp*, *object release* [11]. This discrepancy in definitions presents a challenge when comparing the performance of different approaches.

At the **goal level**, the system observes demonstrations with the aim of identifying the ultimate objectives of tasks rather than the exact procedures used. The key is to infer the reason behind actions, allowing the robot to modify its approach based on its capabilities or the environmental conditions. Zeng et al. proposed an approach to enable robots to understand and execute tasks by interpreting the goals from demonstrations provided by the users [21]. The initial and goal scene are represented by RGB-D images. A method called Discriminatively-Informed Generative Estimation of Scenes and Transforms (DIGEST) was developed to generate scene graph from demonstrated images. The first step in the DIGEST method involves detecting bounding boxes in the RGB images. Once objects are detected, their pose is estimated using the depth information from the RGB-D images and existing object mesh models. The final step is to generate a scene graph. This graph represents the spatial and relational structure of the scene. It is built using inter-object relations such as *exist*, *clear*, *in*, and *on* by calculating the object poses. Given the observation of the goal state of the world, the robot estimates the goal scene graph, and stores the desired inter-object relations by PDDL [22]. It is a formal language used for expressing planning problems and domains. The task planner gives a sequence of high-level pick & place actions. To pick an object, the robot receives a number of pre-computed grasping positions for the object, and uses Moveit!¹ to determine which of these positions can generate a collision-free path.

In conclusion, visual imitation learning leverages the latest advancements in computer vision and machine learning to enhance robot programming at multiple levels. These levels include the execution of basic skills, the understanding of complex task sequences, and the achievement of task goals. For multi-step pick & place tasks, this work implements the imitation learning at the task level, focusing on the integration and sequencing of skills to perform complex, multi-step tasks. Focusing solely on the goal level risks losing important procedural information, while learning at the skill level alone is insufficient for handling multi-step tasks. As introduced in the first section, this framework is realized by mapping the hand trajectories to the robot's end-effector. Comparing to the approach of Wen et al. [18], tracking the hand is also necessary as small objects may be occluded and not visible. In this work, both hand and object trajectories are employed to generate a high-level task plan and low-level trajectories.

¹ <https://moveit.ros.org/>

3. Materials and Methods

3.1. Task Model Representation

To address the three variations mentioned above, a task model is defined by modifying the approach of Backhaus, which was developed for adaptive programming system for assembly tasks [23]. The modified task model consists of four layers, as illustrated in Figure 1. It listed from high-level task plan to low-level position. The task plan is a sequence of actions to complete the task. Each action has a parameter of an object. A *pick* action refers to the act of selecting and lifting an object from a particular location. It involves grasping the object to obtain control over it. The *pick* action is further divided into *Reach* and *Grasp* skills. *Reach* is represented by motion from random start to grasp position. *Grasp* is considered as a discrete event, which occurs in single timestamp. It represents the moment when robot's end-effector successfully makes contact with the object and securely holds it. The *place* action is divided into *Move* and *Release* skills. *Move* is represented by motion from grasp to release positions with the moving object. *Release* is also considered as a discrete event. It involves opening the gripper from the end effector. The skills are extracted from the demonstration video. The grasp and release position of object can be estimated from visual information. This structured task model helps the robot to effectively understand and execute the tasks in dynamic environments.

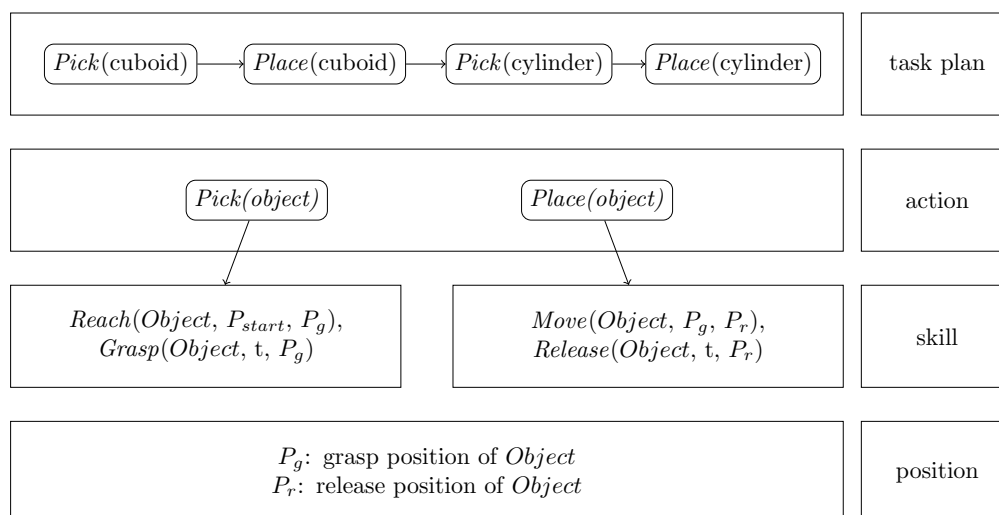


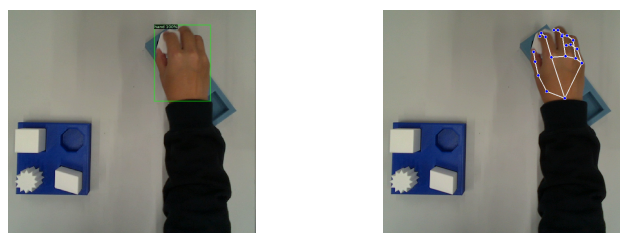
Figure 1. Task model.

3.2. Hand Detection

There are several trained models and frameworks available for hand detection. In this work, two common approaches used in the context of visual imitation learning, one from the Bbox-based approaches and one from the keypoints-based approach are selected to evaluate hand movements. As discussed in the second section, Xin et al. applied a pre-trained Bbox-based model to learn hand motion trajectories with promising results [9]. The model was proposed by Shan et al. to understand hand and hand-object interaction [13]. The system was built on top of a popular object detection system, Faster-RCNN [24]. The model has a two-stage architecture, where the first stage proposes regions and the second stage classifies and refines the bounding boxes. The models were trained on the datasets proposed in the same publication, the 100DOH dataset and the 100K frames. The 100DOH is a large video dataset containing hands and hand-object interactions. They are 27.3K Youtube videos from 11 categories with nearly 131 days of footage of everyday interaction. To detect hand on single image, the authors created a new 100K frame-level dataset, which consists of 99,899 frames extracted from the videos in 100DOH and VLOG Dataset [25]. In total, there are 189,426 annotated bounding boxes for hand. Two models trained on 100K and 100K+ego are provided by the authors for hand

detection. The 100K+ego consists of 56.4K frame subset of egocentric data from [10,26,27]. Both models achieve approximately 90% average precision.

Another detection method is based on keypoints, which is an open-source framework developed by Google called MediaPipe [14]. It provides a pipeline for hand detection. The output of the hand detector are 21 3D hand-knuckle coordinates inside of the detected hand regions. The representation of each keypoint is composed of x -, y - and z - coordinate, where x and y are pixel coordinates on image and z is the distance to the wrist as origin. A hand trajectory can be created from either method using the center of the hand position. The examples for both methods are illustrated in Figure 2a and 2b. The two methods will be evaluated in Section 4.



(a) Detected bounding box using model trained on 100K+ego [13] (b) Detected hand key points with MediaPipe

Figure 2. Results of Hand Detection on RGB Images Using Different Models.

3.3. Object Detection

As discussed in the second section, object detection in visual imitation learning is crucial for understanding task environments. Thanks to the availability of large datasets, various model architectures have been developed and shown to continuously improvement performance. Most of available datasets are labeled with regular bounding boxes for everyday objects. They serve as benchmarks for comparing different model architectures. Some well known datasets are ImageNet [28], Pascal VOC [29] and COCO [30]. They provide raw images, its annotations and standardized evaluation protocols. COCO is currently the standard benchmarking dataset for the object detection community. The latest release of COCO 2017 consists of a total of 123,287 images and 896,782 objects, across its validation datasets. Additionally, the test dataset comprises 40,670 images. They cover 80 object categories, which include everyday objects such as person, bicycle, and car, as well as household items like chair, couch, and dining table. The annotation include the bounding box coordinates, the category level, and other optional attributes. More detailed information about the COCO dataset can be found on the official website². The bounding boxes provided in COCO are rectangular boxes aligned with the image exes that enclose the object. They do not account for the object's rotation. However, the orientation is essential for application scenarios such as understanding aerial images. Dataset of Object deTecton in Aerial images (DOTA) [31] is a large-scale dataset for the application of aerial images. The DOTA dataset consists of a total of 1,793,658 object instances that are annotated with oriented bounding box (OBB) annotations. The illustration of the annotations is shown in Figure 3, where θ refers to the rotation of the rectangular box with respect the image axis. The object instances in DOTA belong to 18 different categories, which include objects like planes, ships, and basketball courts. The OBB approach is followed in this work, which allows for improved accuracy and better object differentiation.

² <https://cocodataset.org>

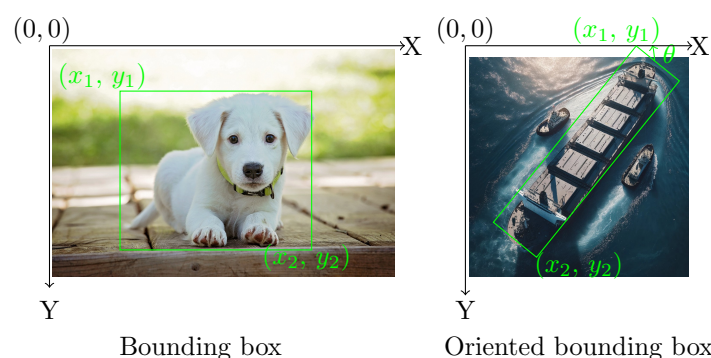


Figure 3. Illustration of the differences between axis-aligned Bounding Box (Bbox) and Oriented Bounding Box (OBB) annotation.

As discussed above, existing datasets primarily cover everyday objects. To develop object detection models for a specific task in manufacturing scenarios, it is necessary to generate a specific set of objects. The data generation process involves capturing images and manually annotating the objects of interest, which is time-consuming. CAD files are usually available in manufacturing scenarios. Synthetic training data can be generated using CAD data and has gained popularity in recent years due to its advantages in providing accurately labeled datasets at a lower cost [32]. In this work, an approach using rendered images in Blender to generate the dataset is used. The details of the synthetic data generation are discussed below.

3.3.1. Data Generation

Photorealistic images are created or rendered to closely resemble real-life objects or scenes. Hodañ et al. developed an approach to synthesize photorealistic images of 3D models, which are used to train a convolutional neural network for detecting the objects in real images [33]. The approach is implemented in BlenderProc, which is an open source pipeline to render the images [15]. The pipeline is a Blender extension with Python API with various examples. The image synthesis approach³ for generating LineMOD dataset in BOP challenge is followed in this work. In this approach, objects are arranged inside an empty room. A random photorealistic material from the CC0 Textures⁴ library is assigned to the walls of the room, and light with a random strength and color is emitted from the room ceiling and from a randomly positioned point light source. Realistic object poses are achieved by dropping objects on the ground plane using PyBullet physics engine integrated in Blender. In the pipeline, both the images and their corresponding annotations are produced simultaneously and automatically. The annotations are generated in OBB format to preserve the orientation of the objects. Examples of photorealistic images of 6 CAD models are shown in Figure 4, which are grayBox, blueBox, cuboid, parallelogram, star and octagon.



Figure 4. Generated photorealistic images.

³ https://github.com/DLR-RM/BlenderProc/blob/main/examples/datasets/bop_challenge/README.md

⁴ <https://ambientcg.com/>

YOLOv8 is utilized in this work to detect object of interests due to the availability of pre-trained models for OBB detection and its software framework. The model utilizes a convolutional neural network that can be divided into two main parts: backbone and detection head. The head of YOLOv8 consists of multiple convolutional layers followed by a series of fully connected layers. These layers are responsible for predicting the oriented bounding boxes and class probabilities for the objects detected in the image. During fine-tuning, the parameters from the model backbone are used to initialize the model, the detection head are initialized with number of classes for the new task. The fine-tuning of the YOLOv8 model with the photorealistic images and the validation on real images will be discussed in the fourth section.

3.4. Mapping Motion from a Human Hand to a Robot End-Effector

This section focuses on mapping motion from a human hand to a robot end-effector. Figure 5 illustrates the mapping concept. The framework starts by generating hand and object trajectories from continuous frame-by-frame hand and object detection. These trajectories are then segmented based on the interaction of the hand and object and the object state. The resulting hand motion segments *reach* and *move* are represented as skills with DMPs. These representations are stored in the task model to reproduce the task.

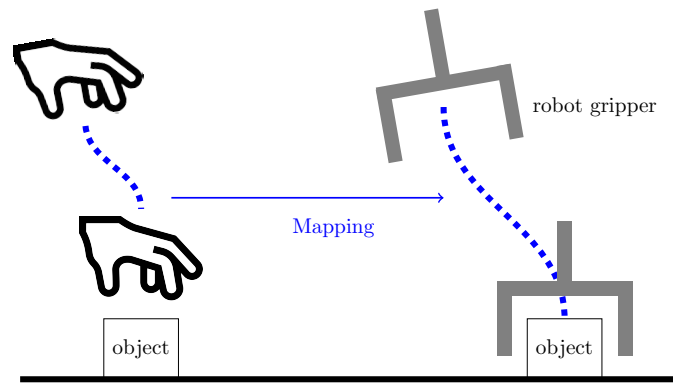


Figure 5. Mapping motion from a human hand to a robot end-effector.

3.4.1. Trajectory Generation for the Human Hand and the Objects

As mentioned in the first section, an RGB-D camera is required to record the visual demonstration. Hand and object detection are performed on the RGB frames. The detected results are bounding boxes at pixel coordinates. The corresponding depth values can be read from the depth frame. The Equation (1) shows the process. The center of detected bounding boxes from RGB images for the human hand and objects are X_{center} and Y_{center} . The responding depth value Z_{center} can be identified from depth image. Similar as RGB image, the depth image is represented by a two-dimensional matrix, $\mathcal{I} = (y, x)$, the coordinate (x, y) corresponds to the column and row indices, respectively. The depth value Z_{center} can be identified from depth image $\mathcal{I}(Y_{center}, X_{center})$. The resulted sequence of $\mathcal{P} = \{X_{center}, Y_{center}, Z_{center}\}$ are hand and object trajectories in pixel coordinates as shown in Figure 6. In the next step, the position are transformed to real world coordinates with camera as origin. In the equation, f_x and f_y describe the focal length of the image, and the ppx and ppy describe the pixel coordinates of the principal point. They are camera intrinsic parameters to describe the internal characteristics of the camera.

$$\begin{aligned} X &= \frac{X_{center} - ppx}{f_x \cdot Z_{center}} \\ Y &= \frac{Y_{center} - ppy}{f_y \cdot Z_{center}} \\ Z &= Z_{center} \end{aligned} \quad (1)$$

According to Ghidoni [34], several technologies are available for estimating depth information: stereoscopy, infrared light and LiDAR. A stereoscopic camera captures two images simultaneously from different angles to determine distance. Infrared and LiDAR technologies work on a similar principle. They use a light source and a receiver to measure the time it takes the light to reach the sensor. Both of these technologies are referred to as ToF sensors. Each technique can be affected by lighting conditions and the texture of objects, and have invalid depth values. To overcome bad depth value, an additional estimated depth map from the color image is used in this work. It uses the global-local path networks model for monocular depth estimation [35]. The model is made available by Hugging Face⁵. In case the depth value is unavailable in aligned depth frame from camera, it will be estimated from color frame using the pre-trained model from Hugging Face. The transformation from pixel coordinates to camera frame is also performed using Equation (1).

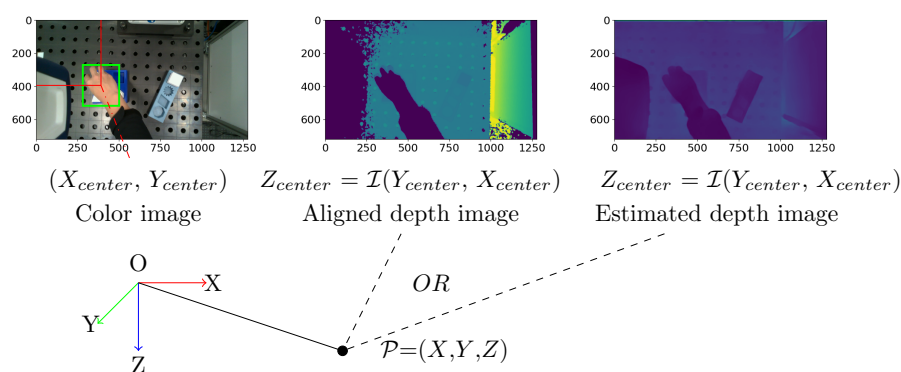


Figure 6. Transform hand and object position from pixel coordinates to camera frame in world coordinates.

By repeating the above mentioned process for each pair of RGB and depth frames, the hand and object trajectories can be generated in 3D space. Normally, the object detection accuracy cannot reach 100%. If the hand and the object of interest cannot be detected in some frames, the last valid detection is propagated to the next frame. The generated hand and object trajectories can be represented as $Traj = (\mathcal{P}_1, \dots, \mathcal{P}_i, \dots, \mathcal{P}_T)$.

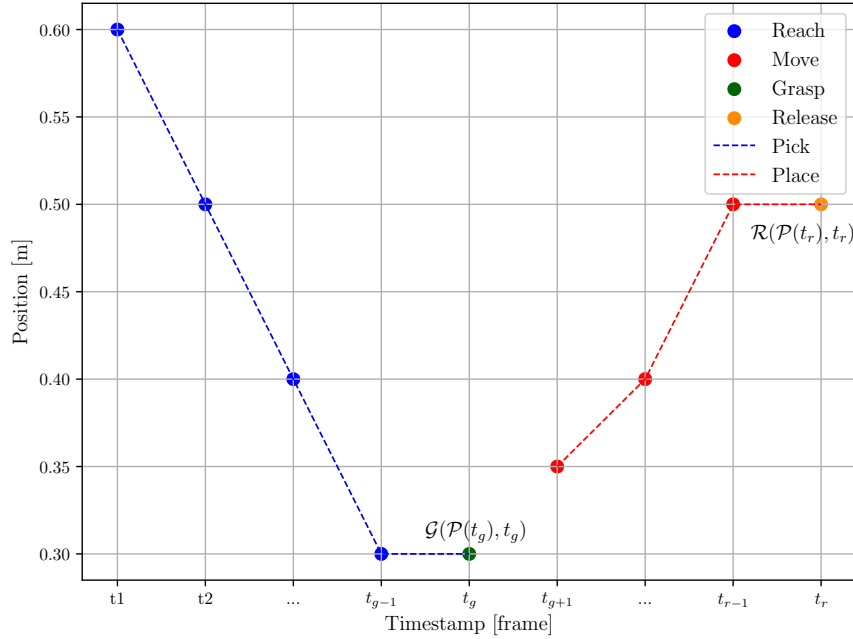
3.4.2. Motion Representation for a Robot's End-Effector

To map hand motions to the robot's end-effector, the hand trajectories are segmented into smaller components which allows adapting hand to end-effector position. The required small components are *pick* and *place* actions and the *reach*, *grasp*, *move* and *release* skills, as defined in Section 3.1. The skills representation for the robot's end-effector during the pick-and-place process is summarized in Table 1, where $\mathcal{P}(t)$ denotes the position of the end-effector at timestamp t . The action of picking up an object involves the *reach* and *grasp* skills. The *reach* skill describes the motion of the end-effector from a starting point to the grasping position. Once the end-effector reaches the object, the *grasp* occurs at the next timestamp. After grasping the object, the *move* skill is used to transport the object to the desired location. Unlike the separate motions of moving and positioning in human hands, transporting an object with a robot is considered as continuous motion. Upon arriving at the desired position, the *release* skill is triggered to release control of the object. An example of the segmented trajectories are illustrated in Figure 7.

⁵ <https://huggingface.co/vinvin02/glpn-nyu>

Table 1. Representation of motion as skills for robots.

Actions	Skills	Segmented hand trajectories
Pick	<i>Reach</i>	$[\mathcal{P}(t_1), \mathcal{P}(t_2), \dots, \mathcal{P}(t_{g-1})]$
	<i>Grasp</i>	$\mathcal{G}(\mathcal{P}(t_g), t_g), \mathcal{P}(t_{g-1}) = \mathcal{P}(t_g)$
Place	<i>Move</i>	$[\mathcal{P}(t_{g+1}), \mathcal{P}(t_{g+2}), \dots, \mathcal{P}(t_{r-1})]$
	<i>Release</i>	$\mathcal{R}(\mathcal{P}(t_r), t_r), \mathcal{P}(t_{r-1}) = \mathcal{P}(t_r)$

**Figure 7.** Illustration of the segmented trajectories.

In the practice of robot programming, after defining the grasp and release positions, additional pre- and post-motion segments are created to ensure a controlled grasp and release. The pre-grasp involves approaching the object at reduced speed and corrected orientation. Once the robot reaches the pre-grasp pose, it performs the *grasp* to securely hold the object. After successfully grasping the object, the robot moves back to post-grasp pose with the reversed approach trajectory to avoid collision with the environment. The pre- and post-grasp pose are typically the same, ensuring consistency in the robot's pose. Once in the pose-grasp pose, the robot transitions to the *move* skill to transport the object to the desired place position. The pre-release involves preparing the robot and the object for release. It includes adjusting the robot's pose to ensure a controlled release. The speed of pre-release and pose-release is usually reduced compared to the transport motion. After releasing the object, the robot returns to post-release pose to start the next movements. Similar to the pre- and pose-grasp poses, the pre- and post-release poses are typically the same for consistency, as shown in Figure 8. In summary, to map a human hand trajectory to a robot's end effector, the hand trajectories are first segmented in to *reach*, *grasp*, *move* and *release*, and then augmented with *pre-grasp*, *post-grasp*, *pre-release* and *post-release*. The orientation $\mathcal{O}(t)$ of the objects is determined during the *pre-grasp* and *pre-release* phases and remains constant throughout the *post-grasp* and *post-release* phases.

$$\begin{aligned}
\mathcal{P}_{pre-grasp} &= \mathcal{P}_{post-grasp} = \mathcal{P}(t_g) \pm \Delta p \\
\mathcal{O}_{pre-grasp} &= \mathcal{O}_{post-grasp} = \mathcal{O}(\text{object}, t_g) \\
\mathcal{P}_{pre-release} &= \mathcal{P}_{post-release} = \mathcal{P}(t_r) \pm \Delta p \\
\mathcal{O}_{pre-release} &= \mathcal{O}_{post-release} = \mathcal{O}(\text{object}, t_r)
\end{aligned}$$

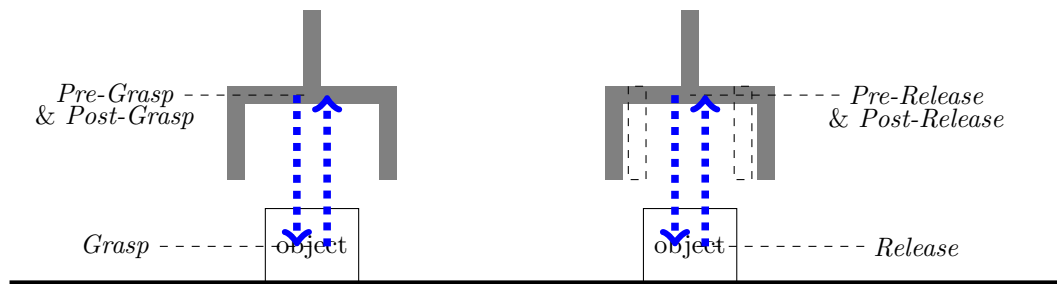


Figure 8. Illustration of Pre-Grasp, Post-Grasp, Pre-Release, Post-Release Phases.

3.4.3. Trajectory Segmentation

This subsection explains how hand trajectories are segmented into smaller elements to represent the skills of *Reach*, *Grasp*, *Move*, and *Release*. This is accomplished by identifying the timestamps for grasping and releasing. The input data consists of the hand and object trajectories generated in Section 3.4.1. Segmenting hand trajectories directly can be a challenging task. Because hand movements often transition smoothly from one phase to another without clear and distinct boundaries. It is therefore difficult to identify where one segment ends and another begins. To address this issue, an approach is proposed that leverages the status of hand-object interaction and object status.

As discussed in Section 3.4.1, before segmenting the trajectories, missing values caused by detection failure are filled by propagating the last valid value to the next one. Then, a Kalman filter is applied to reduce noise. The hand trajectories are first segmented into pick and place cycles for each object by identifying the release timestamps. They are achieved by segmenting the object trajectories by Gaussian Mixture Model (GMM). It is a statistical model that represents a probability distribution $p(x, y)$ between input x and output y variables as a weighted sum of Gaussian distributions.

$$p(\mathcal{X}, y | \theta) = \sum_{k=1}^K p(w_k) p(\mathcal{X}, y | \mu_k, \Sigma_k) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathcal{X}, y | \mu_k, \Sigma_k), \quad (2)$$

where each Gaussian component $\mathcal{N}(x, y | \mu_k, \Sigma_k)$ has prior probability $p(w_k) = \pi_k$, mean μ_k , and covariance matrix Σ_k . GMM parameters $\theta = \pi_k, \mu_k, \Sigma_k$ are learned from training data using the expectation-maximization. It can be utilized for motion segmentation by modeling the distribution of trajectories over time. In a pick & place task, the object trajectory can be divided into three distinct parts, the static state before picking, the moving state, and the static state after placing. The first timestamp in the state after placing indicating the release of the hand. In this work, the object trajectories are segmented by the GMM model with $K = 3$ and $\mathcal{X} = \text{Traj} = \{\mathcal{P}_1, \dots, \mathcal{P}_i, \dots, \mathcal{P}_T\}$.

In the second step, the *grasp* timestamps are used to segment the pick and place motion into *reach* and *move*. When a human hand and other objects are present in the same visual frame, occlusion can occur, which means the hand might partially or fully block the view of the object. This can result in inaccurate object trajectories that do not accurately represent the grasping process. However, the object can be correctly detected at the beginning and end of the demonstration, where the human demonstration has finished or not yet started. The *grasp* timestamp is then determined by calculating the shortest distance between the hand and the object's initial position. The results are *grasp* and

release index in the time sequence for each object. The resulted sequence symbols for three objects can be represented as $[Reach(O_1), Grasp(O_1), Move(O_1), Release(O_1), Reach(O_2), Grasp(O_2), Move(O_2), Release(O_2), Reach(O_3), Grasp(O_3), Move(O_3), Release(O_3)]$.

The segmentation approach is summarized in Algorithm 1.

Algorithm 1 Segmentation hand trajectories

Input: $objectSequence = [3,5,2,4]$ \triangleright The sequence of object labels moved by the robot
Input: $handTraj = (p_1, p_2, \dots, p_T), p_i = (x_i, y_i, z_i), i \in (1, T)$
Input: $objTraj = (o_1, o_2, \dots, o_T), o_i = (x_i, y_i, z_i), i \in (1, T)$
Output: $\{objectLabel: [graspIndex, releaseIndex]\}$ \triangleright A dictionary for all input objects
 $handTraj \leftarrow smoothTraj(handTraj)$ \triangleright The trajectory is processed by kalman filter
 $objTraj \leftarrow smoothTraj(objTraj)$ \triangleright The trajectory is processed by kalman filter
 $start, end \leftarrow process(handTraj)$ \triangleright The first and the last frame with detected hand
for $i = 1$ to T **do**
 $gmm \leftarrow GaussianMixture(n_components = 3)$
 $gmm.fit(objTraj)$
 $labels \leftarrow gmm.predict(objTraj)$
 $releaseIndex \leftarrow$ the first frame in the last cluster
 $initialPosition \leftarrow mean(objTraj[: start])$
 $finalPosition \leftarrow mean(objTraj[end :])$
 $dists \leftarrow$ list
 for x in $hand_traj[start : releaseIndex]$ **do**
 $dist \leftarrow$ euclidean distance between x and $initialPosition$
 Append $dist$ to $dists$
 end for
 $minValue \leftarrow$ minimum value in $dists$
 $minIndex \leftarrow$ index of $minValue$ in $dists$
 $graspIndex \leftarrow minIndex + start$
 $start \leftarrow releaseIndex$
end for

3.4.4. Trajectory Learning & Generation

DMP can learn from a single demonstration and adapt to new start and goal positions [11]. They are therefore used to model goal-directed behaviors. The resulting *reach* and *move* segments are represented by DMPs, which are formalized as stable nonlinear attractor systems [20]. There are many variations of DMPs. As summarized by Fabisch [36], they have in common that

- they have an internal time variable (phase), which is defined by a so-called canonical system,
- they can be adapted by tuning the weights of a forcing term and
- a transformation system generates goal-directed accelerations.

The formulation of DMPs can be found in [16,36,37]. The choice of coordinates for representing a model can make a big difference in how the generalization of the dynamics systems appears [16]. The DMPs are represented in the camera frame in both the trajectory learning and generation processes. Figure 9 illustrates the steps for generating robot end-effector trajectories using DMP. To generate the trajectory for *reach*, the process begins with defining the current position of the end-effector and aims to reach the pre-grasp position. The learned DMPs from the hand trajectories are updated based on defined initial and goal position of the target object. This process generates a smooth trajectory for the specified duration. Similarly, DMP for *move* is updated by defining the post-grasp position and the pre-release position. The adaptability of DMP allows for updating with new initial and goal position, ensuring controlled movements in dynamic environment.

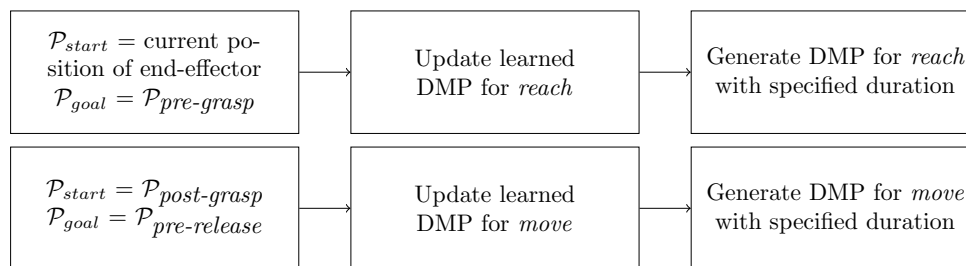


Figure 9. Trajectory generation process.

4. Experimental Results

4.1. Experimental Setup

The goal of this work is to extract the required information from a single video. To avoid any potential bias or singular occurrences, five sets of instructions and their corresponding videos were recorded and evaluated. Figure 10 (a) illustrates the experimental setup to record the videos. The Intel® RealSense™ L515 3D camera is mounted on the robot to record the demonstration process. The camera is positioned to capture the task environment and hand movements. The task environment is shown in Figure 10 (b). As a LiDAR camera, it projects an infrared laser at 860 nm wavelength as an active light source. 3D data is obtained by evaluating the time required for the projected signal to bounce off the objects in the scene and come back to the camera [38]. The frame rate is 30 frames per second (FPS). The videos were recorded using Intel RealSense Viewer and saved as bag files, which can be extracted as sequence of color and depth images. The size of color images recorded by the L515 is 1280×720 and the size of depth image is 640×480 . The depth frame is re-scaled and aligned⁶ to color frame, such that the depth value can be read by pixel coordinate of color image.

The videos feature the same four pick & place tasks, which involve transporting a cuboid, star, parallelogram, and octagon from a blue box to a gray box. At the beginning and end of the recording, the task environment was observed without the presence of hands, ensuring a clear view of the working environment for accurate analysis of the conditions before and after the manipulation was performed. The task was performed with only one hand.



Figure 10. (a) Demonstration setup; (b) Color image; (c) Aligned depth image.

4.2. Evaluation of Hand Detection Methods

In Section 3.2, two methods for hand detection are discussed. The first method is the bounding box-based approach developed by Shan et al. [13], which was trained with the Faster-RCNN model. The second method is the key points-based approach from Google, known as MediaPipe. Both methods were applied to the RGB image sequences. After detection, the results were further processed by

⁶ <https://dev.intelrealsense.com/docs/rs-align-advanced>

transforming the hand locations into Cartesian coordinates in the camera’s coordinate system. This transformation leverages the depth value of the center point of the detected hand.

The models provided by Shan et al. [13] are available online⁷ as open source. The detection process is integrated into Detectron2⁸, which is a framework for computer vision tasks such as object detection and instance segmentation. The model trained on 100K + ego was used to run the test. The score threshold for filtering out low-confidence detection was set at 0.7. Any model output below this threshold is considered as invalid. A low threshold means that the model accepts detections with lower confidence levels, which can significantly increase the number of false positives. The performance of the both methods are evaluated on the five recorded videos, which consists of 3601 images in total, where hand are visible in 2631 images. The results are summarized in Table 2. The percentages in the table represent the proportion of images where the methods successfully detected hands when hands are present. It was observed that the model failed when the hand was too small in the image, an example is shown in Figure 11.



Figure 11. An instance of missed hand detection by the model Faster-RCNN.

The key points-base method MediaPipe⁹ is an open-source framework, which is available for use with different programming languages and devices. The hand landmarks detection solution in Python within this framework is used to evaluate the performance. The parameter “static image mode” was set to True and the “maximum number of hands” was set to 1. The results in Table 2 show that the Faster-RCNN is generally more reliable for hand detection tasks than MediaPipe.

Table 2. Comparison of hand detection accuracy on recorded videos between Faster-RCNN and Mediapipe.

	Faster-RCNN	MediaPipe
Video 1	100.00%	93.99%
Video 2	95.76%	8.86%
Video 3	99.82%	15.74%
Video 4	92.59%	9.94%
Video 5	91.67%	12.50%

The hand trajectories are then computed using the detection results from Faster-RCNN because it has better performance. The approach was presented in Section 3.4.1. An example of a hand trajectory is shown in Figure 12. From this graph, the hand movement is most variable along the z-axis (depth), as indicated by the green line. This is due to the hand moving towards and away from the camera for the multi-step pick & place task.

⁷ https://github.com/ddshan/hand_detector.d2
⁸ <https://ai.meta.com/tools/detectron2/>
⁹ https://developers.google.com/mediapipe/solutions/vision/hand_landmarker

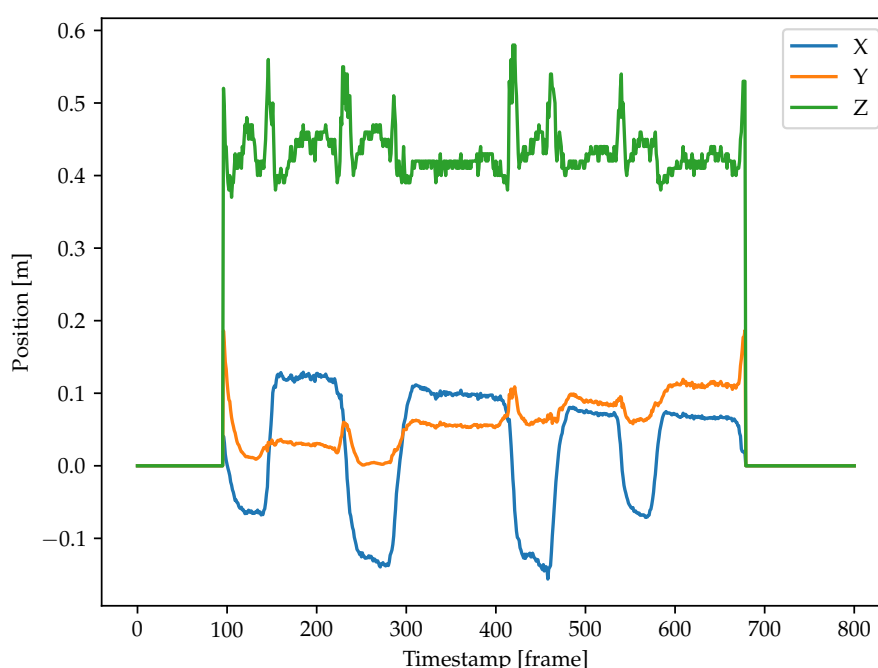


Figure 12. Generated hand trajectory in camera frame with Faster-RCNN.

4.3. Training and Evaluation the Object Detection Model

When it comes to detecting objects that are not as commonly modeled as hands, it is typically necessary to generate a custom dataset for training or fine-tuning. A total of 18,750 photorealistic images were generated using the approach discussed in Section 3.3.1 to train the YOLOv8 model for the OBB detection task. Out of these, 15,000 images are used for the training set, and 3,750 images are allocated for the validation set to evaluate the performance of the model. Transfer learning approach is used with the pre-trained model YOLOv8 from Ultralytics¹⁰. During training, the number of epochs was set to 20. The IoU was set as 0.7. The learning rate was set to the default value of 0.01, as this has been demonstrated to be a value for transfer learning tasks.

The Table 3 presents the object detection results for the validation set consisting of 3,750 images with a total of 21,821 instances across various classes. The overall precision and recall for detecting all classes are 98.5% and 99% respectively, with a mAP at 50% IoU threshold of 98.7% and mAP at 50-95% IoU of 47.8%. Each class, including GrayBox, BlueBox, Parallelogram, Cuboid, Octagon, and Star, has a high precision and recall, ranging from 96.8% to 99.3% and 97.4% to 99.4% respectively. The mAP50 for these classes varies slightly, between 98.4% and 99%, while the mAP50-95 ranges from 43.4% to 49.6%. The Star class has the highest mAP50-95 at 49.6% and shares the highest precision and recall with the Octagon class at 99.3%. The Cuboid class shows a slightly lower precision compared to others at 96.8%. Overall, the results indicate a high level of accuracy in the object detection task for this validation set.

¹⁰ <https://docs.ultralytics.com/tasks/obb/#models>

Table 3. Object detection results on validation set.

Classes	Images	Instances	Precision	Recall	mAP50	mAP50-95
all	3750	21821	98.5%	99%	98.7%	47.8%
GrayBox	3750	3633	98.8%	99.1%	98.6%	43.4%
BlueBox	3750	3670	98.7%	99.3%	99%	48.4%
Parallelogram	3750	3637	98.3%	97.4%	98.4%	48.3%
Cuboid	3750	3598	96.8%	99.4%	98.4%	48.4%
Octagon	3750	3652	99.1%	99.3%	98.7%	48.6%
Star	3750	3631	99.3%	99.3%	99%	49.6%

The object detection model trained on photorealistic images was evaluated on real-world images to assess its performance. Three representative images with the detection results are presented in Figure 13. In these images, the model is able to identify and localize the object of interest within the scene. However, the occlusion and the presence of a human hand show poor performance.

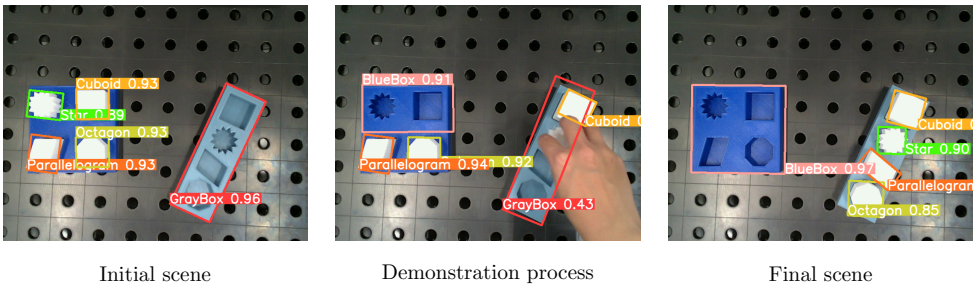


Figure 13. Object detection results on real world images with Faster-RCNN.

To evaluate the performance of the domain adaption, images without the presence of a human hand were considered. 971 images from the 5 recorded videos are extracted. The Table 4 presents performance for different classes, including the number of images, instances and accuracy. The model achieves impressive performance on the entire dataset, with 100% accuracy indicating that all detected instances are correct. The recall of 83.3% suggests that the model is capturing a significant fraction of the actual instances. The mAP50 score of 90.5% indicates a high average precision at a threshold of 50. The model performs differently for each class. For instance, the GrayBox and BlueBox class exhibit lower recall values of approximately 50% due to object occlusion within the box. It is worth noting that the mAP50-95 scores are relatively lower compared to the mAP50 scores, indicating potential for improvement in detecting instances within this threshold range.

Table 4. Object detection results on real world images.

Classes	Images	Instances	Precision	Recall	mAP50	mAP50-95
all	971	5826	100%	83.3%	90.5%	72.1%
GrayBox	971	971	100%	50.3%	75.1%	38.5%
BlueBox	971	971	99.9%	49.7%	70.1%	70.1%
Parallelogram	971	971	100%	100%	99.5%	73.5%
Cuboid	971	971	100%	100%	99.5%	71.4%
Octagon	971	971	100%	100%	99.5%	89.9%
Star	971	971	100%	100%	99.5%	89.3%

4.4. Evaluation of Depth Value Quality

Depth value is essential to generate the hand and the objects trajectories in Cartesian coordinates. Chen et al. discusses the common challenge of motion blur in depth sensing technology, particularly in dynamic movements with ToF sensors [39]. Depth is calculated by measuring the time it takes for a light signal to travel to the object and back to the sensor. However, movement of objects can distort the timing measurement, resulting in inaccuracies. The experimental results of calculating the depth value of hand and object show the same pattern, where the depth value of hand trajectories is generally worse than that of object trajectories. Specifically, the study notes that valid depth values for hand trajectories fall below 50%, whereas for objects, they exceed 70%. In case of invalid depth value, the depth is estimated using the method discussed in Section 3.4.1. The difference between the depth value estimated by the pre-trained deep learning model and the valid depth value from the camera is less than 1cm. This evaluation demonstrates the necessity of using depth estimation techniques to ensure accurate trajectory generation and highlights the performance of the applied depth estimation method.

4.5. Evaluation of the Proposed Segmentation Approach

The hand trajectories from the recorded videos were then segmented using the Algorithm 1. They are calculated by the object trajectories and the distance between the hand and the target object. An object trajectory is outlined in Figure 14. Each color represents a component generated by GMM. The trajectories were taken as correct segmented by inspecting the image at the identified timestamps for *grasp* and *release*. This method of verification ensures that the segmentation aligns with the actual moments when the object is manipulated. The approach was successful in identifying the *grasp* and *release* events for all target objects in the video with high accuracy.

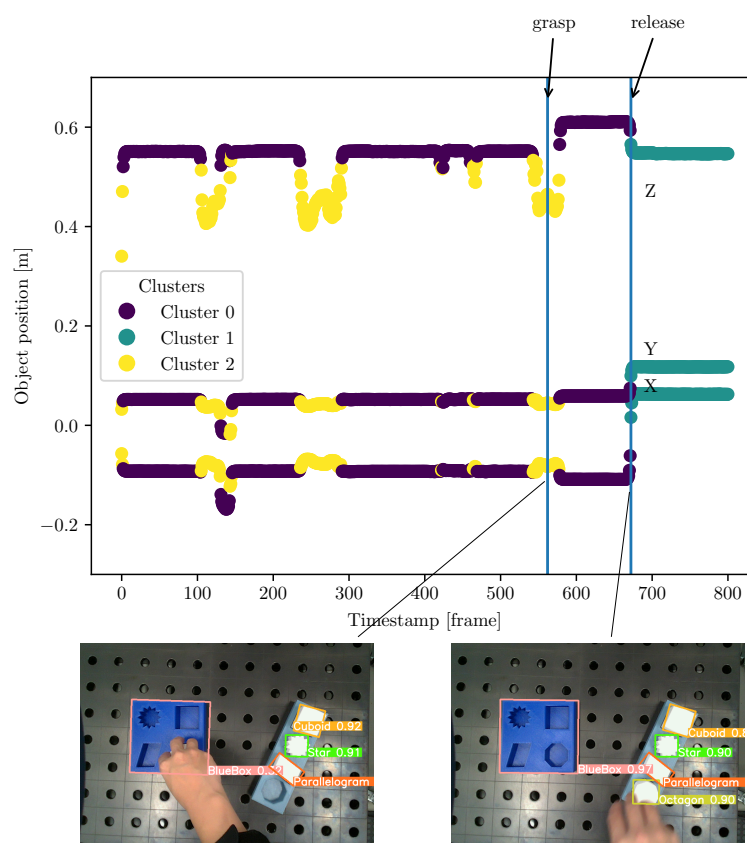


Figure 14. Hand trajectory segmentation results.

4.6. Evaluation of the Motion Learning Approach

As discussed in Section 3.4.2, the segmented trajectories are modeled by DMP, which represent *reach* and *move* skills. As shown in Figure 12, the hand trajectories have irregularities, which illustrates the necessity for DMP to smooth these trajectories. Moreover, the DMP should be capable of generalizing to new start and target positions. The implementation of Fabisch [36] in Python was used to learn the DMPs¹¹. Figure 15 illustrates the learned and updated DMP for representing the skill *reach* with regulation coefficient $\lambda = 0$. The number of ϕ functions was set to $N = 20$. The “Demonstration” curve represents a segmented trajectory of a hand motion from the starting point to a *pre-grasp* position. The “Reproduction” curve shows how the DMP was learned to imitate the original hand motion by attempting to follow the demonstrated trajectory. To enhance the manipulation accuracy, the DMP was adapted to object’s position “New Start” and “New Goal”. Since $\lambda = 0$, the noise and irregularities in the hand trajectory are encoded in the DMP, the phase approaching the goal in the figure shows the description. Figure 16 illustrates the improved DMP with $\lambda = 0.1$, the “Reproduction” and “Adaptation” curves show that the influence of noise from the demonstration trajectory is reduced. Figure 17 shows the learned and updated DMP for representing the *move*. In summary, the learned DMP are able to adapt learned behaviors to new situations while preserving the details of the original demonstration.

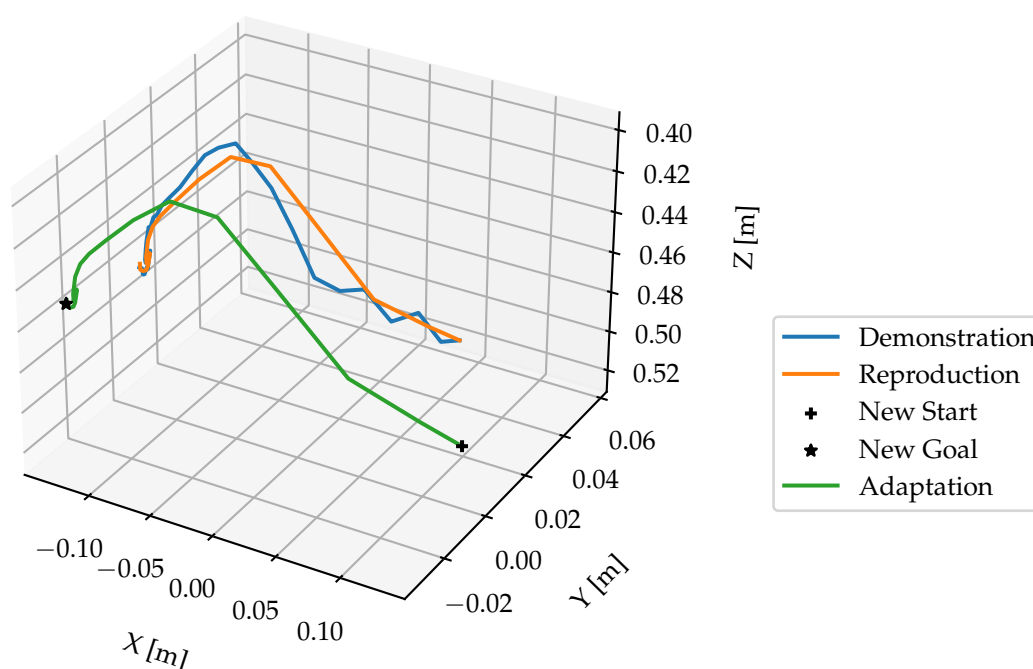


Figure 15. Learned and updated DMP for representing *reach* with $\lambda = 0$ and $N = 20$

¹¹ https://github.com/dfki-ric/movement_primitives

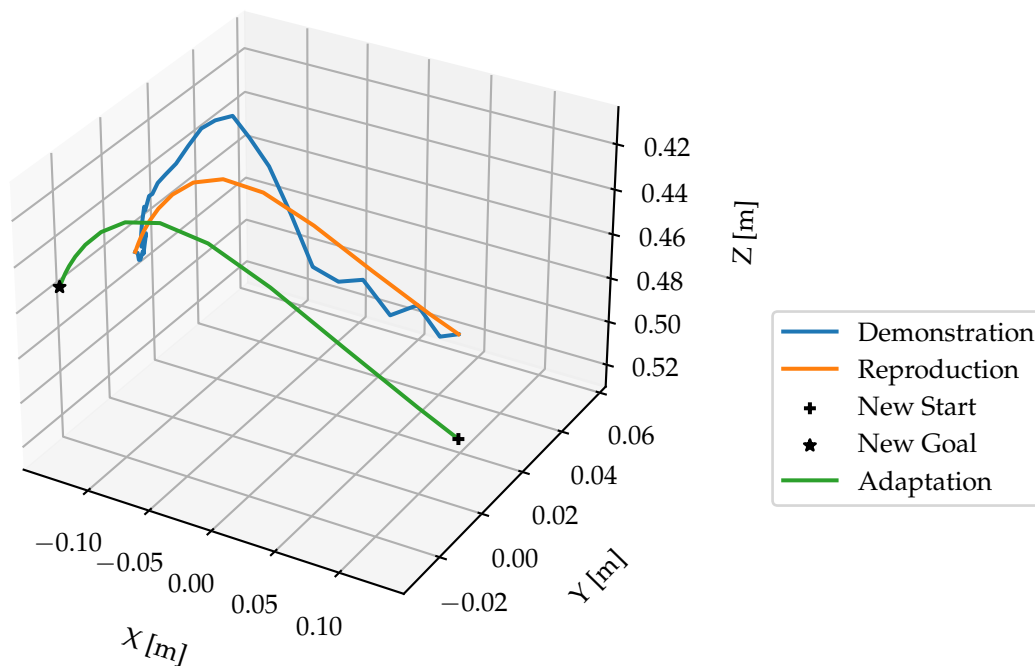


Figure 16. Learned and updated DMP for representing *reach* with $\lambda = 0.1$ and $N = 20$

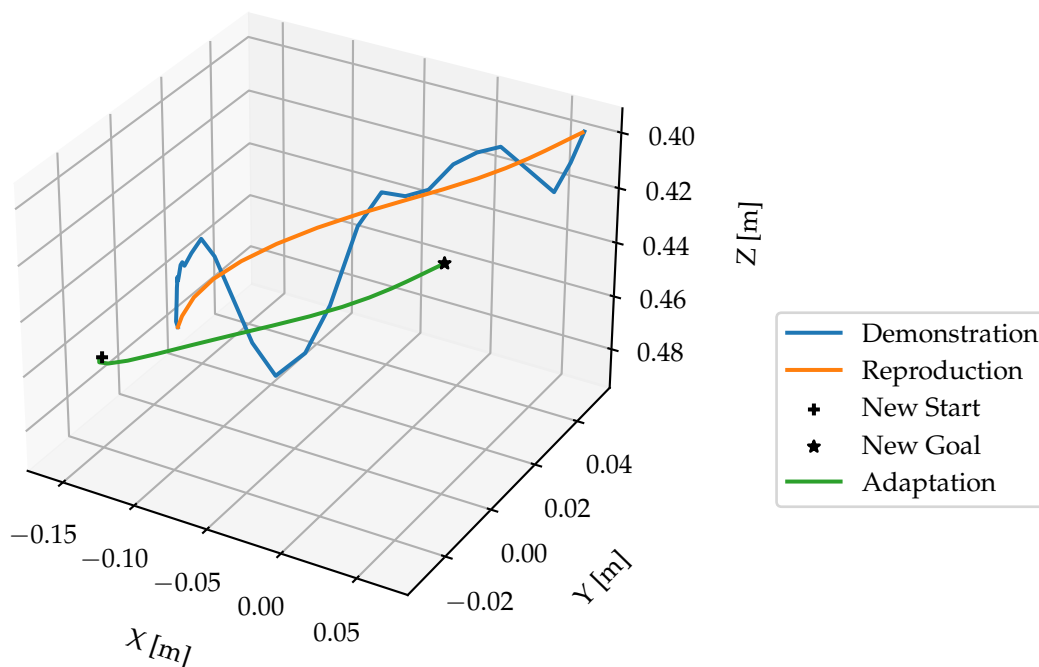


Figure 17. Learned and updated DMP for representing *move* with $\lambda = 0.1$ and $N = 20$.

4.7. Computational Time

The evaluation was conducted on a computer running Ubuntu 20.04, equipped with an AMD Ryzen 5 3600 6-core processor (12 threads), 31.3 GiB of RAM, and a GeForce RTX 2070 Super GPU. Table 5 provides a breakdown of tasks organized into three phases: Preparation, Learning,

and Execution. The time figures presented are averaged from the five recorded videos. In the preparation phase, tasks such as training data generation and fine-tuning the YOLOv8 model for object detection takes the most time, 10 hours and 2 hours, respectively. Although these tasks are time-consuming, this phase can be automated and performed independent of task learning. The learning phase starts with video recording, which takes 30 seconds, followed by converting the video to color and depth frames (70 seconds), and depth estimation, which is the most time-consuming task in the phase at 1.2 hours. After that, the system performs hand detection (18 seconds) and object detection (20 seconds). It then generates hand and object trajectories, both taking just 1 second each. This is followed by trajectory segmentation, which requires 4 seconds, and finally, the system learns dynamic movement primitives (DMPs) in less than a second. Despite the longer duration of depth estimation, the learning phase overall is highly efficient. During execution, the required trajectories can be generated in less than a second.

Table 5. Task breakdown with one video which contains of 597 frames.

Phase	Task	Time
Preparation	Training data generation	10 hours
	Fine-tuning the YOLOv8 model	2 hours
Learning	Video recording	30 seconds
	Video to color and depth frame	70 seconds
	Depth estimation	1.2 hour
	Hand detection	18 seconds
	Object detection	20 seconds
	Hand trajectory generation	1 second
	Object trajectory generation	1 second
	Trajectory segmentation	4 seconds
	Learning DMP	< 1 second
Execution	Generation DMP	< 1 second

4.8. Evaluation on Robot in Simulation

The evaluation results show a 100% success rate, with all five recorded tasks successfully executed within the ROS environment (see Figure 18). The grasp and release positions for each object were defined by the simulation environment and the DMPs were generated accordingly. This performance highlights the effectiveness of the one-shot visual imitation learning framework, which demonstrates robust learning of human actions. In addition, the system shows adaptability to changes in object positions, further reinforcing its effectiveness in dynamic environments.

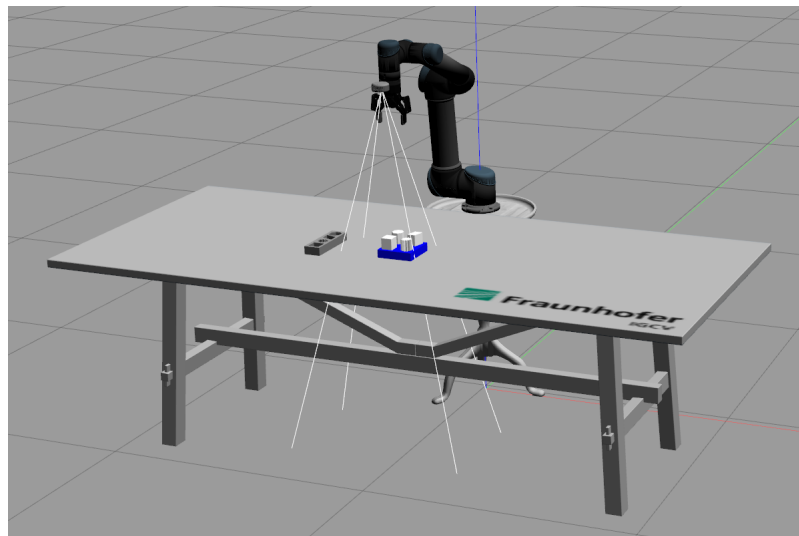


Figure 18. Simulation environment.

5. Discussion

The results demonstrate the potential of one-shot visual imitation learning for reducing the complexity and time required for robot programming in industrial settings. The modular design of the framework allows for flexibility, enabling individual components to be improved without affecting the overall system. However, several challenges remain, particularly in real-world object detection and trajectory optimization.

Faster-RCNN proved to be reliable for hand detection in controlled scenarios. However, limitations were observed when the hands were small from the camera, highlighting the need for human movements with hands clearly visible during demonstrations. For object detection, the use of synthetic data generated from CAD models was successful, but domain adaptation remains crucial for achieving higher accuracy in real-world tasks. The proposed segmentation algorithm effectively identified grasp and release events, demonstrating its effectiveness in task decomposition. DMPs provided a robust foundation for skill learning, allowing the robot to adapt to new task configurations. However, optimizing the generated trajectories for efficiency remains an area for future improvement. Additionally, faster depth estimation techniques is required to improve the computational efficiency.

6. Conclusions

This paper presents a one-shot visual imitation learning framework designed to enable robots to learn multi-step pick-and-place tasks from a single video demonstration. The key contributions include a modular approach that integrates hand detection, object detection, trajectory segmentation, and skill learning through DMPs. The framework demonstrates significant potential in reducing the data requirements and simplifying the robot programming process, making it particularly valuable for manufacturing applications where flexibility and efficiency are essential.

Overall, the framework has demonstrated its ability to learn from one-shot demonstrations and provides a practical solution for intuitive robot programming in manufacturing contexts. Future work should focus on improving the robustness of object detection, enhancing trajectory optimization, and reducing the computation time for depth estimation to increase the system's overall effectiveness.

In conclusion, the proposed framework is a significant step toward more flexible and efficient robot programming by non-experts. Its modularity ensures that individual components can be further optimized and adapted to suit a wide range of applications.

Author Contributions: Conceptualization, Shuang Lu; methodology, Shuang Lu; software, Shuang Lu; validation, Shuang Lu; writing—original draft preparation, Shuang Lu; writing—review and editing, Christian Härdtlein,

Johanness Schilp; visualization, Shuang Lu; supervision, Johanness Schilp; All authors have read and agreed to the published version of the manuscript.

Funding: This research was conducted within the project MeMoRob (AKZ: DIK0358/01) funded by the Bavarian Ministry of Economic Affairs, Regional Development and Energy (StMWi). The authors would like to thank the StMWi for the financial support.

Informed Consent Statement: Not applicable

Data Availability Statement: The data is available at <https://drive.google.com/drive/folders/11LpzuVkabRRSv2kR7raaqmaqvvZ52s05?usp=sharing>.

Acknowledgments: We would like to extend our gratitude to our colleagues at Fraunhofer IGCV, Christian Hädrlein and Julian Glaß, for providing the CAD files and the objects used in the experiments.

Conflicts of Interest: The authors declare no conflicts of interest.

References

1. Beibl, J.; Lee, J.; Krause, D.; Moon, S.K. Flexibility - Grand Challenge for Product Design and Production: Review and Status. *Procedia CIRP* **2023**, *119*, 91–96. The 33rd CIRP Design Conference, doi:https://doi.org/10.1016/j.procir.2023.05.003.
2. Steinmetz, F.; Nitsch, V.; Stulp, F. Intuitive task-level programming by demonstration through semantic skill recognition. *IEEE Robot. Autom. Lett.* **2019**, *4*, 3742–3749.
3. Berg, J.K. System zur aufgabenorientierten Programmierung für die Mensch-Roboter-Kooperation. PhD thesis, Technische Universität München, 2020.
4. Saveriano, M. Robotic Tasks Acquisition via Human Guidance: Representation, Learning and Execution. PhD thesis, Technische Universität München, 2017.
5. Qiu, Z.; Eiband, T.; Li, S.; Lee, D. Hand Pose-based Task Learning from Visual Observations with Semantic Skill Extraction. 2020 29th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN). IEEE, 2020, pp. 596–603.
6. Finn, C.; Yu, T.; Zhang, T.; Abbeel, P.; Levine, S. One-shot visual imitation learning via meta-learning. Conference on Robot Learning. PMLR, 2017, pp. 357–368.
7. Nagarajan, T.; Feichtenhofer, C.; Grauman, K. Grounded human-object interaction hotspots from video. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 8688–8697.
8. Zeng, Z. Semantic Robot Programming for Taskable Goal-Directed Manipulation. PhD thesis, University of Michigan, 2020.
9. Xin, J.; Wang, L.; Xu, K.; Yang, C.; Yin, B. Learning Interaction Regions and Motion Trajectories Simultaneously from Egocentric Demonstration Videos. *IEEE Robotics and Automation Letters* **2023**.
10. Damen, D.; Doughty, H.; Farinella, G.M.; Fidler, S.; Furnari, A.; Kazakos, E.; Moltisanti, D.; Munro, J.; Perrett, T.; Price, W.; Wray, M. Scaling Egocentric Vision: The EPIC-KITCHENS Dataset. European Conference on Computer Vision (ECCV), 2018.
11. Kyrarini, M.; Haseeb, M.A.; Ristić-Durrant, D.; Gräser, A. Robot learning of industrial assembly task via human demonstrations. *Auton. Robot.* **2019**, *43*, 239–257.
12. Ding, G.; Liu, Y.; Zang, X.; Zhang, X.; Liu, G.; Zhao, J. A Task-Learning Strategy for Robotic Assembly Tasks from Human Demonstrations. *Sensors* **2020**, *20*, 5505.
13. Shan, D.; Geng, J.; Shu, M.; Fouhey, D.F. Understanding human hands in contact at internet scale. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 9869–9878.
14. Lugaresi, C.; Tang, J.; Nash, H.; McClanahan, C.; Uboweja, E.; Hays, M.; Zhang, F.; Chang, C.L.; Yong, M.G.; Lee, J.; others. Mediapipe: A framework for building perception pipelines. *arXiv preprint arXiv:1906.08172* **2019**.
15. Denninger, M.; Winkelbauer, D.; Sundermeyer, M.; Boerdijk, W.; Knauer, M.; Strobl, K.H.; Humt, M.; Triebel, R. BlenderProc2: A Procedural Pipeline for Photorealistic Rendering. *J. Open Source Softw.* **2023**, *8*, 4901, doi:10.21105/joss.04901.
16. Ijspeert, A.J.; Nakanishi, J.; Hoffmann, H.; Pastor, P.; Schaal, S. Dynamical movement primitives: Learning attractor models for motor behaviors. *Neural Comput.* **2013**, *25*, 328–373.
17. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster r-cnn: Towards real-time object detection with region proposal networks. *Adv. Neural Inf. Process. Syst.* **2015**, *28*.

18. Wen, B.; Lian, W.; Bekris, K.; Schaal, S. You Only Demonstrate Once: Category-Level Manipulation from Single Visual Demonstration. *Robotics: Science and Systems 2022* **2022**.
19. Li, S.; Lee, D. Point-to-pose voting based hand pose estimation using residual permutation equivariant layer. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11927–11936.
20. Schaal, S. Dynamic movement primitives-a framework for motor control in humans and humanoid robotics. In *Adaptive motion of animals and machines*; Springer; pp. 261–280.
21. Zeng, Z.; Zhou, Z.; Sui, Z.; Jenkins, O.C. Semantic robot programming for goal-directed manipulation in cluttered scenes. *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 7462–7469.
22. Aeronautiques, C.; Howe, A.; Knoblock, C.; McDermott, I.D.; Ram, A.; Veloso, M.; Weld, D.; Sri, D.W.; Barrett, A.; Christianson, D.; others. Pddl the planning domain definition language. *Technical Report, Tech. Rep.* **1998**.
23. Backhaus, J.C.S. Adaptierbares aufgabenorientiertes Programmiersystem für Montagesysteme. PhD thesis, Technische Universität München, 2016.
24. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster r-cnn: Towards real-time object detection with region proposal networks. *Adv. Neural Inf. Process. Syst.* **2015**, 28.
25. Fouhey, D.F.; Kuo, W.; Efros, A.A.; Malik, J. From Lifestyle VLOGs to Everyday Interactions. *CVPR*, 2018.
26. Li, Y.; Ye, Z.; Rehg, J.M. Delving into egocentric actions. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 287–295.
27. Sigurdsson, G.A.; Gupta, A.K.; Schmid, C.; Farhadi, A.; Karteek, A. Actor and Observer: Joint Modeling of First and Third-Person Videos. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* **2018**, pp. 7396–7404.
28. Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; others. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vis.* **2015**, 115, 211–252.
29. Everingham, M.; Van Gool, L.; Williams, C.K.; Winn, J.; Zisserman, A. The pascal visual object classes (voc) challenge. *Int. J. Comput. Vis.* **2010**, 88, 303–338.
30. Lin, T.Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; Zitnick, C.L. Microsoft coco: Common objects in context. *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V* 13. Springer, 2014, pp. 740–755.
31. Ding, J.; Xue, N.; Xia, G.S.; Bai, X.; Yang, W.; Yang, M.; Belongie, S.; Luo, J.; Datcu, M.; Pelillo, M.; Zhang, L. Object Detection in Aerial Images: A Large-Scale Benchmark and Challenges. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2021**, pp. 1–1, doi:10.1109/TPAMI.2021.3117983.
32. Vanherle, B.; Moonen, S.; Reeth, F.V.; Michiels, N. Analysis of Training Object Detection Models with Synthetic Data. *33rd British Machine Vision Conference 2022, BMVC 2022, London, UK, November 21–24, 2022*. BMVA Press, 2022.
33. Hodaň, T.; Vineet, V.; Gal, R.; Shalev, E.; Hanzelka, J.; Connell, T.; Urbina, P.; Sinha, S.N.; Guenter, B. Photorealistic image synthesis for object instance detection. *2019 IEEE international conference on image processing (ICIP)*. IEEE, 2019, pp. 66–70.
34. Ghidoni, S. Performance evaluation of depth completion neural networks for various RGB-D camera technologies **2023**.
35. Kim, D.; Ga, W.; Ahn, P.; Joo, D.; Chun, S.; Kim, J. Global-Local Path Networks for Monocular Depth Estimation with Vertical CutDepth. *CoRR* **2022**, abs/2201.07436, [2201.07436].
36. Fabisch, A. Learning and generalizing behaviors for robots from human demonstration. PhD thesis, Universität Bremen, 2020.
37. Pastor, P.; Hoffmann, H.; Asfour, T.; Schaal, S. Learning and generalization of motor skills by learning from demonstration. *2009 IEEE International Conference on Robotics and Automation*. IEEE, 2009, pp. 763–768.

38. Servi, M.; Mussi, E.; Profili, A.; Furferi, R.; Volpe, Y.; Governi, L.; Buonamici, F. Metrological Characterization and Comparison of D415, D455, L515 RealSense Devices in the Close Range. *Sensors* **2021**, *21*, 7770.
39. Chen, Z.; Liu, P.; Wen, F.; Wang, J.; Ying, R. Restoration of Motion Blur in Time-of-Flight Depth Image Using Data Alignment. 2020 International Conference on 3D Vision (3DV), 2020, pp. 820–828, doi:10.1109/3DV50981.2020.00092.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.