



INNOVATIVE ROBOTICS FOR AGILE PRODUCTION

D3.1: Planner guided by demonstration extracted from a monocular camera



PROJECT ACRONYM: AGIMUS

PROGRAMME: Horizon Europe

GRANT AGREEMENT: No 101070165

TYPE OF ACTION: Horizon Research & Innovation Actions

START DATE: 1 October 2022

DURATION: 48 months



Funded by
the European Union

Document information

Issued by:	CTU
Issue date:	01/09/2023
Due date:	30/09/2023
Work package leader:	CNRS
Start date:	1 October 2022
Dissemination level:	PUB (Public)

Document History

Version	Date	Modifications made by
1.0	01/09/2023	1st draft version distributed - CTU

Authors

First name	Last name	Beneficiary
Vladimir	Petrik	CTU
Kateryna	Zorina	CTU
Josef	Sivic	CTU
Florent	Lamiroux	CNRS

In case you want any additional information, or you want to consult with the authors of this document, please send your inquiries to: vladimir.petrik@cvut.cz.

Quality reviewers

First name	Last name	Beneficiary
Justin	Carpentier	INRIA
Maximilien	Naveau	TOWARD

Disclaimer

Funded by the European Union under GA no. 101070165. Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the European Union. The granting authority cannot be held responsible for them.

© AGIMUS Consortium, 2022

Reproduction is authorised provided the source is acknowledged.

Executive summary

This report was developed in the context of the AGIMUS project, funded by the Horizon Europe Framework Program for Research and Innovation of the European Union for 2021-2027. Its objective is to present the progress mainly towards the specific objective **SpO1.1: Learn complex multi-step manipulation tasks from demonstrations of human experts recorded using low-cost sensors** that is part of the **Objective 1: Significantly accelerate the introduction and deployment of a robotic system to a new agile production environment through advanced perception and understanding of human-centric feedback.**

Table of Contents

1 Introduction	6
2 Demonstration Extraction from Video	6
3 Planner Guided by Demonstration	7
4 Conclusion	8
Reference List	8

List of Figures

1 Approach overview.	6
2 Configuration space	7

1 Introduction

The presented work is based on our paper entitled **Multi-Contact Task and Motion Planning Guided by Video Demonstration** [3], published to the 2023 IEEE International Conference on Robotics and Automation (ICRA 2023). The full version of the paper is attached at the end of the deliverable. In the paper [3], we estimated the sequence of 6D poses of known objects from a monocular video as presented in more detail in Sec. 2. This sequence of poses is used to build the sequence of tasks the planner should investigate as discussed in Sec. 3. The overview of the proposed approach is shown in Fig. 1. Construction of the admissible configuration space from the video demonstration guides the planner towards the poses used in the demonstration and results in faster motion planning. In addition, video guidance allowed us to solve task-and-motion-planning problems that were not solvable with state-of-the-art planners, as we demonstrated on a newly created challenging benchmark. We made both the benchmark as well as planner publicly available at https://github.com/agimus-project/guided_tamp_benchmark and https://github.com/agimus-project/video_guided_tamp_planner.

2 Demonstration Extraction from Video

The first step of the proposed pipeline (Fig. 1 (i)) extracts the 6D poses of objects from the video recorded by a calibrated monocular camera. The pipeline uses an AI-based hand/object contact detector [2]. The detector is evaluated for each frame of the video to detect when a human grasps or releases an object. At these changes in the contact state, the 6D pose of the object is extracted using an AI-based iterative method called CosyPose [1]. The 3D geometry models of the objects are assumed to be known for CosyPose, which is a reasonable assumption in many practical robotics setups, such as in manufacturing.

This step contributes to the goals of Task **T3.1 Extracting demonstration from monocular cameras** by integrating software for contact state and object pose detections. Estimating

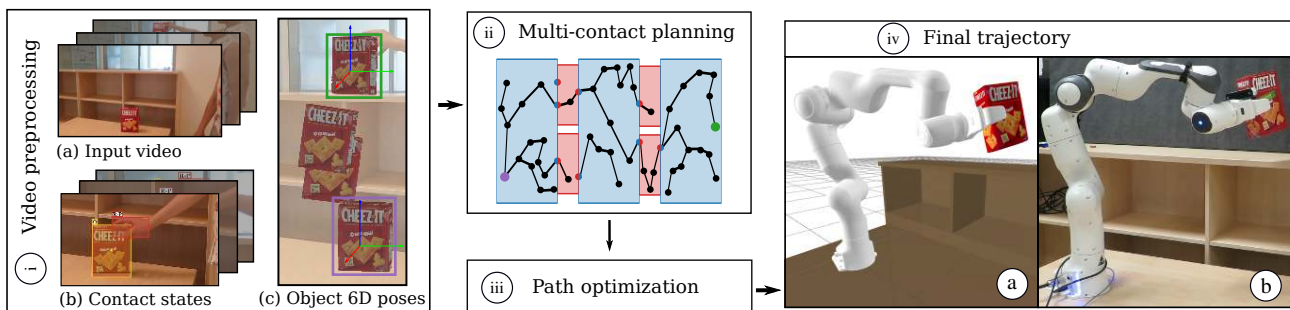


Figure 1: Approach overview. (i) First, we extract contact states and 6D object poses from the input instructional video. (ii) These poses are used to construct admissible configurations, and we grow multiple trees in this space until we find a path between the start and goal configurations. (iii) This path is then further shortened by an optimization module and (iv) executed either in simulation (iv-a) or on a real-world robot (iv-b).

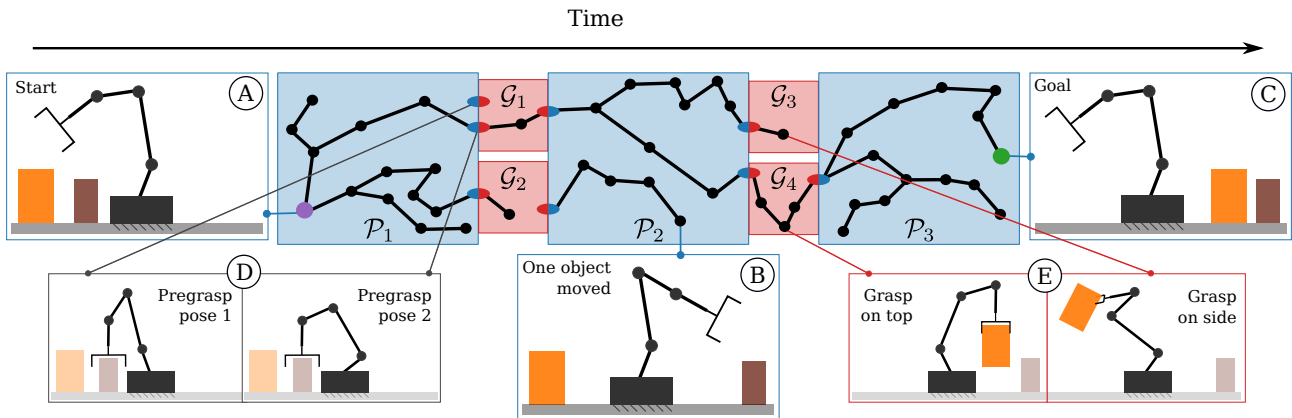


Figure 2: The admissible configuration space is a set of (i) placement states \mathcal{P}_i where objects rest on the contact surfaces at specific poses and (ii) grasp states \mathcal{G}_j where one of the objects is grasped by the robot gripper. To transit to another placement state, *i.e.* to change the pose of one object, the robot needs to grasp the object and release it at the new location. Since we define several ways to grasp each object, we can transition through multiple grasp states to build a path between configurations given by two consecutive placement states. Start configuration (A) lies in the first placement state \mathcal{P}_1 . To achieve the goal, the robot needs to move the brown object first, *i.e.* to reach state \mathcal{P}_2 *e.g.* configuration (B)). Finally, the robot moves the second object (reach \mathcal{P}_3) and moves the robot configuration to the given goal (C). The configuration at the transition between placement state and grasp state is not unique as there are various robot joint values resulting in the same 'pregrasp' poses as shown in (D). Multiple ways of grasping the same object are represented by different grasp states as shown in (E).

human pose and optimizing through the contacts will be implemented in the next version to improve the accuracy of the object pose detection and to provide additional information for the planner.

3 Planner Guided by Demonstration

We use the poses of the object at the contact state changes to construct the configuration space for the planner. The configuration space consists of a sequence of manifolds; each manifold is defined by a set of constraints that respect the information estimated from the video. The example of the configuration space for a planar robot is visualized in Fig. 2. There are two main groups of manifolds: placement manifolds where objects are static on contact surfaces, and grasp manifolds where one of the objects is grasped (attached) to a robot and other objects are static on contact surfaces. To guide the planner, we modified the RRT algorithm to sample multiple trees with roots at the intersection of the manifolds, where root configuration satisfies both grasp and placement constraints at the same time. This helps the planner to discover narrow spaces in the configuration space and therefore, find the solution to challenging planning problems. Our guided planner outperformed state-of-the-art planners that do not use video guidance in the newly created challenging benchmark.

The proposed planner contributes to the goals of **T3.3: Combining demonstration and**

planning by designing a planner that uses poses extracted from the video demonstration. The globally optimal solution is provided by the human depicted in the instructional video and therefore the hard-to-solve combinatorial problem is resolved by the demonstration. As the proposed planner focuses on the motion of objects instead of the human, the motion transfer from human to robot is resolved by design.

4 Conclusion

We present an approach for guiding task and motion planning by a video demonstration. Contact states and 6D object poses extracted from the video are used to build the tree in the configuration space of the planner efficiently. This work makes a step towards large-scale learning of task & motion plans from instructional videos. The software for the planner is open source and available publicly at https://github.com/agimus-project/video_guided_tamp_planner and the benchmark that was used to evaluate its performance is available at https://github.com/agimus-project/guided_tamp_benchmark.

As the proposed planner can solve hard planning problems, it can be used to simplify the transfer of human knowledge in industrial applications, *i.e.* be used to teach the robot. In the context of AGIMUS project, it will be used to generate sample trajectories for task **T3.2: Motion dataset from the planner** and to solve multi-step task-planning problem in **T5.5: Multi-step Task Planning Experiment**.

Reference List

- [1] Yann Labbé et al. “CosyPose: Consistent Multi-view Multi-object 6D Pose Estimation”. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. Ed. by Andrea Vedaldi et al. Cham: Springer International Publishing, 2020, pp. 574–591. ISBN: 978-3-030-58520-4.
- [2] Dandan Shan et al. “Understanding Human Hands in Contact at Internet Scale”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2020.
- [3] Kateryna Zorina et al. “Multi-contact task and motion planning guided by video demonstration”. In: *ICRA 2023-International Conference on Robotics and Automation*. 2023.

Multi-Contact Task and Motion Planning Guided by Video Demonstration

Kateryna Zorina[♣] David Kovar[♣] Florent Lamiroux[◇] Nicolas Mansard[◇]
Justin Carpentier[♥] Josef Sivic[♣] Vladimir Petrik[♣]

Abstract—This work aims at leveraging instructional video to guide the solving of complex multi-contact task-and-motion planning tasks in robotics. Towards this goal, we propose an extension of the well-established Rapidly-Exploring Random Tree (RRT) planner, which simultaneously grows multiple trees around grasp and release states extracted from the guiding video. Our key novelty lies in combining contact states, and 3D object poses extracted from the guiding video with a traditional planning algorithm that allows us to solve tasks with sequential dependencies, for example, if an object needs to be placed at a specific location to be grasped later. To demonstrate the benefits of the proposed video-guided planning approach, we design a new benchmark with three challenging tasks: (i) 3D re-arrangement of multiple objects between a table and a shelf, (ii) multi-contact transfer of an object through a tunnel, and (iii) transferring objects using a tray in a similar way a waiter transfers dishes. We demonstrate the effectiveness of our planning algorithm on several robots, including the *Franka Emika Panda* and the *KUKA KMR iiwa*.

I. INTRODUCTION

Traditional robot motion planning algorithms seek a collision-free path from a given starting robot configuration to a given goal robot configuration [1]. Despite the large dimensionality of the configuration space, sampling-based motion planning algorithms [2], [3] have shown to be highly effective for solving complex motion planning problems for robots, ranging from six degrees of freedom (DoF) for industrial manipulators to tens of DoFs for humanoids [4]. Manipulation task-and-motion planning (TAMP) [5] adds an additional complexity to the problem by including movable objects in the state space. This requires the task-and-motion planner to discover the pick-and-place actions that connect the given start and goal robot configurations bringing the manipulated objects from their start poses to their goal poses. The combined need to plan the continuous robot motion together with discrete grasp and release actions makes it very challenging to find a plan for long sequential tasks involving several pick-and-place actions.

To address this challenge, we propose to leverage human demonstrations of the given task provided by an instructional

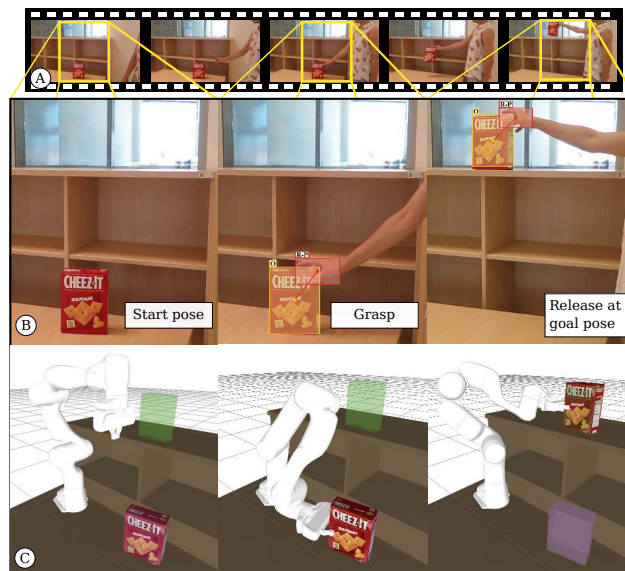


Fig. 1. The proposed planning approach is guided by the demonstration video (A). The video depicts a person manipulating a known object; the *cheez-it* box in this particular example. The video can contain several pick-and-place actions with multiple objects. Here only a short clip with only one object and one action is shown. From the video we recognize (i) the contact states between the human hand and the object, marked by red bounding boxes in (B); and (ii) the object 6D pose (3D translation and 3D rotation w.r.t camera) at the grasp and release contact states, marked in yellow in (B). The robot trajectory planned by the proposed approach is shown in (C). The start and goal object poses in (C) are shown in magenta and green, respectively.

video. Changes of contacts, which are difficult to discover by planning, are extracted from the instructional video together with the 6D poses of *a priori* known objects as visualized in Fig. 1. The extracted information is then used to guide the proposed RRT-like algorithm that simultaneously grows multiple trees around the grasp and release states recognized in the video. To show the benefits of guidance by an instructional video, we design a new benchmark with three challenging multi-contact tasks for manipulation task-and-motion planning that are inspired by real-life problems and are difficult to solve otherwise by state-of-the-art TAMP solvers.

The contributions of this paper are three-fold:

- we propose an RRT-like multi-tree planner that incorporates video demonstration to solve complex multi-contact TAMP tasks;
- we design a new benchmark with three challenging tasks that we will make publicly available to facilitate progress on this challenging problem;
- we provide experimental results demonstrating that our

[♣] CIIRC, Czech Technical University in Prague

[◇] LAAS-CNRS, Université de Toulouse, CNRS, Toulouse

[♥] INRIA, Paris

This work is part of the AGIMUS project, funded by the European Union under GA no.101070165. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Commission. Neither the European Union nor the European Commission can be held responsible for them. This work was funded by the European Regional Development Fund under the project IMPACT (reg. No. CZ.02.1.01/0.0/0.0/15.003/0000468), the Grant Agency of the CTU in Prague, grant No. SGS21/178/OHK3/3T/17, and the Louis Vuitton ENS Chair on Artificial Intelligence..

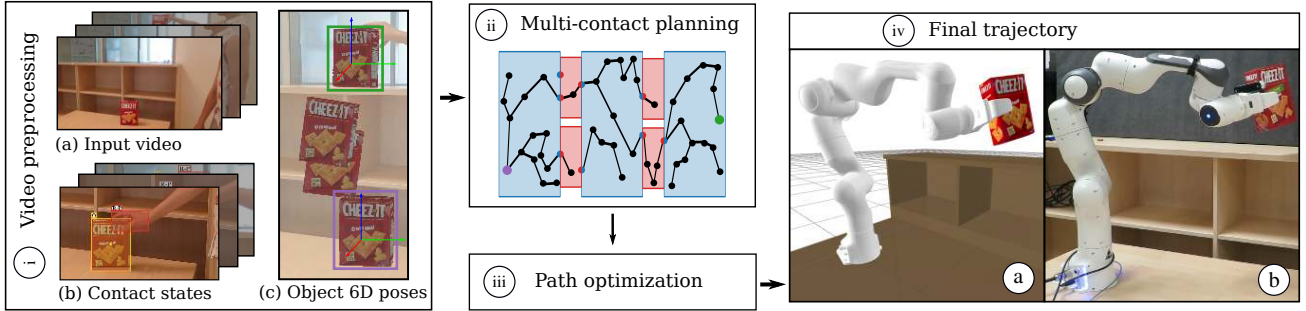


Fig. 2. **Approach overview.** (i) First, we extract contact states and 6D object poses from the input instructional video as described in Sec. III-A. (ii) Next, we grow multiple trees in the admissible configuration space until we find a path between the start and goal configurations. More details on the state space are in Sec. III-B and more details on planning the path in Sec. III-C. (iii) This path is then further shortened by an optimization module, and (iv) executed either in simulation (iv-a) or on a real-world robot (iv-b).

approach outperforms state-of-the-art TAMP solvers on the proposed benchmark.

Please refer to our project page [6] for the supplementary video and code.

II. RELATED WORK

Motion planning. A popular choice for planning the motion of robotic manipulators is to use sampling-based planners inspired by RRT [7], for example, RRT connect [3] that grow trees simultaneously from the start and goal configurations. Our approach is also inspired by RRT but simultaneously grows multiple trees. Using multiple trees for motion planning has been explored in, for example, [8], [9], [10]. In [8], two categories of trees are used: (i) global trees, that are initialized at the start and goal configurations and (ii) local trees, that are initialized from configurations that fail to connect to global trees. Local trees have a higher probability of being constructed in difficult regions such as narrow passages. The trees grow independently until merged with each other. Authors of [9] and [10] explore multi-tree approaches for solving multi-goal planning problems by constructing a tree for every target location. In our approach, the additional trees are sampled at contact state changes that are hard to explore by the planner alone. In our case, the input video demonstration provides the time of the contact state change and the 6D poses of objects for which we sample the ‘local’ trees. This provides a significant advantage compared to the other sampling methods.

Manipulation task-and-motion planning seeks the robot motion that manipulates objects to achieve their given goal poses. This manipulation is often referred to as object rearrangement planning and was studied by several TAMP solvers [5], [11], [12], [13], [14] with some approaches limiting their applicability to 2D tabletop manipulation [12], [13], [14] where objects are picked and placed only on a given flat surface of a table. Compared to these methods, we solve 3D rearrangement tasks where objects are moved outside the flat surface of the table, for example, onto a shelf. An existing approach not limited to 2D rearrangement planning [11] combines symbolic planning with conditional samplers and is effective for multi-contact and cost-sensitive problems. However, it still suffers from the limitations of sampling-based methods. For example, discovering narrow

passages in the configuration space remains a major open challenge. In the proposed work, the input instructional video guides the planner through these narrow passages in the configuration space. Another general TAMP solver is called *Humanoid Path Planner* (HPP) [5] and uses RRT connect [3] to solve TAMP tasks. The planning in HPP is performed in admissible configuration space that is represented by a set of manifolds (denoted *states* in HPP) defined by placement and/or grasp numerical constraints. The movement between the states is achieved through *transitions* that allow *e.g.* to grasp the object at a specific handle. We build on the HPP planner and construct an admissible configuration space based on the input video demonstration.

Guiding planners by demonstration without using video has been explored in several set-ups. For example, several kinesthetic demonstrations obtained by operators guiding the robot manually are used in [15] together with a sampling-based planner to solve household tasks like transferring sugar with a spoon or cleaning a table. Demonstrations in 3D simulation obtained by manipulating a 22 DoFs glove are used in [16] to provide a sequence of actions and to identify parts of the objects that are grasped. The planner is finding a collision-free path for the robot to repeat the sequence of pick-and-place actions. Instead of obtaining demonstrations by kinesthetic guidance or teleoperation by a glove, we propose to extract demonstrations from a video depicting the human performing the given task.

Related to us, a video of a human demonstration has been used to obtain a sequence of commands for a robot in [17]. It is assumed that a robot controller is available to perform these commands, for example, to grasp or to carry an object, or to pour from the cup. Therefore only the task-space action planning is performed. In our approach, we perform task-and-motion planning and do not rely on given internal controllers to perform the individual manipulation actions. Instead, we rely on a contact recognizer [18] and a 6D pose estimator [19] to extract the relevant guidance for the planner from the video. Such contact and object pose estimators from images have been shown to be robust to various lighting conditions, camera properties, etc. While there are several methods to estimate finger-level grasps of objects [20], [21], robots do not typically grasp objects in the same way as humans and therefore the robotic planner may not benefit

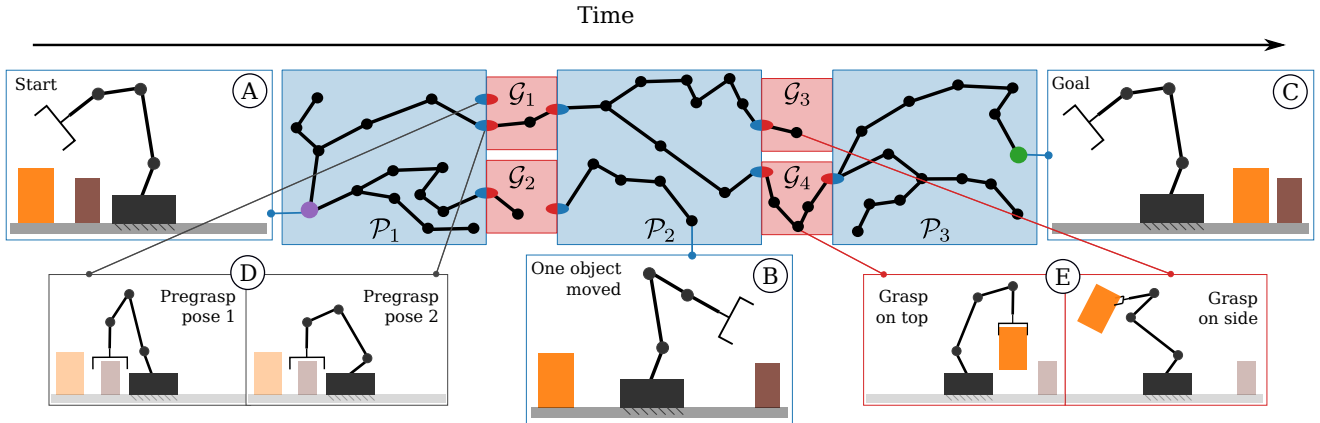


Fig. 3. **The admissible configuration space** is a set of (i) **placement states** \mathcal{P}_i where objects rest on the contact surfaces at specific poses and (ii) **grasp states** \mathcal{G}_j where one of the objects is grasped by the robot gripper. To transit to another **placement state**, *i.e.* to change the pose of one object, the robot needs to grasp the object and release it at the new location. Since we define several ways to grasp each object, we can transition through multiple **grasp states** to build a path between configurations given by two consecutive **placement states**. Start configuration (A) lies in the first **placement state** \mathcal{P}_1 . To achieve the goal, the robot needs to move the brown object first, *i.e.* to reach state \mathcal{P}_2 (*e.g.* configuration (B)). Finally, the robot moves the second object (reach \mathcal{P}_3) and moves robot configuration to the given goal (C). The configuration at the transition between the **placement state** and **grasp state** is not unique as there are various robot joint values resulting in the same 'pregrasp' poses as shown in (D). Multiple ways of grasping the same object are represented by different **grasp states** as shown in (E).

from this information. Hence, we focus only on information related to the pose of objects and the contact states and bypass completely the challenging task of extracting finger-level grasping information from the input video.

Learning-based methods. Policy search methods, for example, reinforcement learning, are widely used to solve manipulation tasks [22], [23], [24]. However, these methods usually require substantial training time, and once the policy is learned, it usually does not generalize well to different scenes, for example, if the furniture is moved. To avoid these limitations, we rely on a geometry-based planning approach that is able to re-plan the path if the geometry of the scene changes.

III. MULTI-CONTACT PLANNING GUIDED BY VIDEO DEMONSTRATION

The proposed method uses a video demonstration to guide the *task & motion* planning algorithm for the given manipulation task. The overall pipeline is visualized in Fig. 2. The input to the proposed approach is the instructional video depicting a human manipulating objects and the geometrical description of the scene. The output of the approach is a collision-free path for the robot that solves demonstrated manipulation task.

A. Extracting contact states and object poses from input video

The objective of this stage is to extract the changes of the contact states, *i.e.* when a human grasps or releases an object, and extract the object 6D poses at the times of the contact state changes. The input video demonstrates a human manipulating known objects one by one from their starting poses toward their goal poses. We focus the proposed method on a single-arm robotic manipulator; therefore, the human demonstrator only moves one object at a time. We apply the hand contact recognizer [18] to obtain a sequence of contact states $c_k \in \{\text{grasp}, \text{release}, \text{none}\}^N$, where N is the number

of movable objects, the contact changes are indexed by $k \in \{1, \dots, T\}$ with T being the number of changes in contact states. Therefore, the variable c_k contains the contact state for each object in the scene at the time of contact state change k . The contact states of the non-manipulated objects are set to 'none'. The 3D models of objects used in the demonstration are assumed to be known, which is a reasonable assumption in many practical robotics setups, such as in manufacturing. As a result, the object 6D poses can be estimated from a single video frame by a render & compare pose estimator such as CosyPose [19]. The 6D pose is estimated for every movable object in the scene at the time of the contact state change k . The poses of the objects are denoted as $A_k \in SE(3)$, where for each of the objects we estimate the 4×4 homogeneous matrix that describes the pose of the object in the common frame of reference. We use the recognized contact states and object poses to constrain the configuration space in which the planning is performed as described next.

B. The configuration space built from contacts and poses

In the case of manipulation TAMP, planning is performed in an admissible configuration space, which is defined by a set of numerical constraints that prevent objects from flying in the air or moving without being grasped by the robot. In [5], the space of admissible configurations is represented by a set of states, *i.e.* a set of manifolds, in which the admissible configurations lie. Here a configuration is defined by the position and orientation of objects in the scene as well as the configuration of the robot. The configurations that belong to these states are subject to placement and/or grasp constraints. The placement constraint for an object enforces that the object remains in the same stable pose on the contact surface. A grasp constraint for an object enforces that the object remains static w.r.t. to the gripper frame. For each object, we define several possible handles, *i.e.* several transformations between the robot gripper and the object. We define *placement states* as those that are

subject to only the placement constraint and *grasp states* as those that are subject to the grasp constraint for one of the objects and placement constraints for the rest of the objects. The transition between the placement and the grasp states represents configurations in which an object is grasped but lies on the contact surface at the same time, *i.e.* constraints from the placement and grasp states are satisfied simultaneously. We call such states *neighboring states*.

To simplify the planning, we propose simplifying the admissible configuration space by using the recognized contact states and object poses extracted from the instructional video. For each unique set of static object poses estimated from the video, we construct a placement state. The constructed placement states represent the sequence of object placements observed in the video. The consecutive placement states are *connected* by grasp states as the object needs to be grasped and moved to change the placement state, *i.e.* to change the pose of one object. There are multiple handles for each object, and therefore there are multiple grasp states, *i.e.* multiple ways to grasp a specific object. An example of such an admissible configuration space is shown in Fig. 3, where three placement states are created: the first placement state \mathcal{P}_1 is constrained (given) by the initial poses of objects, the next placement state \mathcal{P}_2 is constrained by the starting pose of the orange object and the goal pose of the brown object, and the last placement state \mathcal{P}_3 is constrained by the goal poses of both objects. In the example, two handles were defined for each object, representing (i) grasp on top and (ii) grasp on the side. For the path to be feasible, it must contain only the configurations from the admissible space and satisfy additional constraints of the environment, *e.g.* robot configuration must be within the joint limits, and the configuration must be collision-free. We refer to the space in which the configurations satisfy all restrictions as C_{free} .

C. Planning between contact states

Transitioning between the placement and grasp states, which is necessary to move the objects, remains a challenge for sampling-based planners such as RRT [3] as it needs to sample configurations that satisfy constraints from both neighboring states. To address that issue, we design an extension of RRT that grows multiple trees simultaneously with tree roots sampled at the transitions between placement and grasp states. The overview of the proposed algorithm is shown in Alg. 1. The algorithm is split into two main routines: (i) sampling of a new tree at the transition between the placement and grasp states, and (ii) growing an existing tree at a randomly selected state. We randomly select which routines are used in each iteration with a Bernoulli distribution controlled by a parameter $\eta_{\text{sample_tree}}$. The algorithm stops if both start and goal configurations are connected into a single tree or if we run out of resources (*e.g.*, the maximum number of iterations or maximum planning time).

Sample from the state. One of the main capabilities required by the algorithm is sampling from the given state. Similarly to [5], we sample a random configuration from Euclidean

Algorithm 1 Multi-contact RRT guided by demonstration

Require: Contact states c_k , object 6D poses A_k , $\mathbf{q}_{\text{start}}$, \mathbf{q}_{goal} , $\eta_{\text{sample_tree}}$, step size δ

```

1:  $\mathcal{T} \leftarrow \{\text{init\_tree}(\mathbf{q}_{\text{start}}), \text{init\_tree}(\mathbf{q}_{\text{goal}})\}$   $\triangleright$  Existing trees
2: repeat
3:    $p \sim \mathcal{U}(0, 1)$ 
4:   if  $p < \eta_{\text{sample\_tree}}$  then  $\triangleright$  Sampling a new tree
5:      $k \sim \{1, \dots, T\}$ 
6:      $\mathbf{q}_{\text{from}}, \mathbf{q}_{\text{to}} \leftarrow \text{sample\_on\_transition}(c_k, A_k)$ 
7:     if  $\mathbf{q}_{\text{from}} \notin C_{\text{free}}$  or  $\mathbf{q}_{\text{to}} \notin C_{\text{free}}$  then
8:       continue
9:     end if
10:     $t \leftarrow \text{init\_tree}(\mathbf{q}_{\text{from}})$ 
11:     $t.\text{add\_edge}(\mathbf{q}_{\text{from}}, \mathbf{q}_{\text{to}})$ 
12:     $\mathcal{T} \leftarrow \mathcal{T} \cup \{t\}$ 
13:     $\text{attempt\_link}(\mathbf{q}_{\text{from}}, \mathcal{T}, \delta)$   $\triangleright$  Alg. 2
14:     $\text{attempt\_link}(\mathbf{q}_{\text{to}}, \mathcal{T}, \delta)$   $\triangleright$  Alg. 2
15:  else  $\triangleright$  Tree growing
16:     $\mathcal{S} \leftarrow \text{sample\_state}(\mathcal{T})$ 
17:     $\mathbf{q}_{\text{rand}} \leftarrow \text{sample\_configuration}(\mathcal{S})$ 
18:     $\mathbf{q}_{\text{nn}} \leftarrow \text{nearest\_neighbor}(\mathbf{q}_{\text{rand}}, \mathcal{S}, \mathcal{T})$ 
19:     $\mathbf{q}_{\text{step}} \leftarrow \mathbf{q}_{\text{nn}} + \delta(\mathbf{q}_{\text{rand}} - \mathbf{q}_{\text{nn}})$ 
20:    if  $\mathbf{q}_{\text{step}} \notin C_{\text{free}}$  then
21:      continue
22:    end if
23:     $t \leftarrow \text{get\_tree}(\mathbf{q}_{\text{nn}})$ 
24:     $t.\text{add\_edge}(\mathbf{q}_{\text{nn}}, \mathbf{q}_{\text{step}})$ 
25:     $\text{attempt\_link}(\mathbf{q}_{\text{step}}, \mathcal{T}, \delta)$   $\triangleright$  Alg. 2
26:  end if
27: until  $\text{get\_tree}(\mathbf{q}_{\text{start}}) = \text{get\_tree}(\mathbf{q}_{\text{goal}})$  or out of resources

```

space and call a numerical solver [25] iteratively to compute the configuration that satisfies the numerical constraints of the state. To sample the transition connecting two states, the constraints from both states are merged. However, to assign a unique state to each configuration sampled from the transition, we construct two identical configurations denoted \mathbf{q}_{from} and \mathbf{q}_{to} and assign states to them. The two main routines of Alg. 1, which use the described sampling procedure, are discussed next.

Sampling a new tree at the transition between the placement and grasp states is performed as follows. First, we sample a pair of configurations $(\mathbf{q}_{\text{from}}, \mathbf{q}_{\text{to}})$ that satisfy the constraints of the randomly selected transition. If both configurations also satisfy the constraints of the environment (*i.e.* respect the joint limits and are collision-free), we create a new tree containing a root (\mathbf{q}_{from}) and a leaf (\mathbf{q}_{to}). We attempt to link both the created configurations to the existing trees in their corresponding states.

Tree growing in a randomly selected state is performed as follows. This routine consists of sampling a new configuration in admissible configuration space and extending the tree in the direction of the sampled configuration. We randomly choose a state \mathcal{S} and sample the random configuration \mathbf{q}_{rand} from it. The nearest neighbor of the sampled

configuration q_{nn} is found such that q_{nn} also lies in the state \mathcal{S} . A new configuration q_{step} is computed along the segment between q_{nn} and q_{rand} in the manually defined step-size distance δ from q_{nn} . If q_{step} is collision-free, we add it to the tree. Finally, we attempt to link the new configuration to existing trees with nodes in the same state.

Attempt to link function is used by Alg. 1 to connect given configuration q with the existing trees. The function is shown in Alg. 2. It searches for a linear path between the configuration q and other tree nodes in the same state. For each tree, the nearest neighbor q_{nn} is found. Then, configurations along the segment from q_{nn} to q are added to the tree unless a collision is observed. If the entire path is collision-free, we merge the trees that contain the configurations q_{nn} and q .

Algorithm 2 Attempt to link function connects a given configuration to other trees

Require: configuration q , set of trees \mathcal{T} , step size δ ,

```

1:  $\mathcal{S} \leftarrow \text{get\_state}(q)$ 
2: for  $t \in \mathcal{T} \setminus \text{get\_tree}(q)$  do
3:   if  $t$  has nodes in  $\mathcal{S}$  then
4:      $q_{nn} \leftarrow \text{nearest\_neighbor}(q, \mathcal{S}, \{t\})$ 
5:      $q_{parent} \leftarrow q_{nn}$ 
6:     for  $q_{step} \in \{q_{nn}, q_{nn} + \delta, \dots, q_{nn} + n\delta, q\}$  do
7:       if  $q_{step} \notin C_{free}$  then
8:         break
9:       end if
10:       $t.\text{add\_edge}(q_{step}, q_{parent})$ 
11:       $q_{parent} \leftarrow q_{step}$ 
12:    end for
13:    if  $q_{step} = q$  then
14:      merge  $t$  and tree containing  $q$ 
15:    end if
16:  end if
17: end for

```

D. Path optimization

The proposed multi-contact RRT produces a feasible path in the admissible configuration space that achieves the goal. However, this path often contains redundant motion, for example, approaching an object several times before grasping it. To further optimize the path, we apply an approach similar to [26]: we select two random configurations q_1 and q_2 that lie on the path and belong to the same state. We try to build a collision-free segment between them. If such a segment exists, we replace the portion of the path between q_1 and q_2 by this segment.

IV. EXPERIMENTS

In this section we first describe the new benchmark for multi-contact task & motion planning and then give the details of our evaluation including a discussion of the benefits and limitations of our approach.

A. Benchmark for multi-contact task & motion planning

We designed a new benchmark consisting of three challenging tasks that involve multiple contact changes, *i.e.* ob-

ject pick-and-place actions. The objects used for manipulation were selected from the YCBV dataset [27] for which the 6D pose estimator [19] is available. The camera used for the video demonstration recording was calibrated both intrinsically and extrinsically with respect to the environment. This allows us to express object poses and furniture poses in a common frame of reference. Multiple robots with varying number of degrees of freedom are tested in each task, including the Panda robot and the Kuka IIWA robot mounted on a mobile platform. For the purpose of benchmarking, we sample the attachment pose of the non-mobile robots to the desk randomly, such that all pick-and-place poses in the demonstration are reachable. For real scenarios, the robot-to-desk attachment is given by the environment. Next, we provide the description of the tasks. Please see the supplementary video for their visualization.

The shelf task is composed of a table, a shelf and a varying set of objects that the robotic manipulator is supposed to arrange, *i.e.* to move the objects to the predefined poses on the table or on the shelf. The complexity of the task is controlled by the number of objects that the robot should arrange. This task is challenging for state-of-the-art planners because it requires moving multiple objects in a single planning task.

The tunnel task consists of a tunnel and a single object that should be transferred through the tunnel. The tunnel is thin enough so that the robot can place an object inside the tunnel on one side and pick it up from the other side. The challenge lies in the narrow passage in the admissible configuration space that needs to be discovered by the planner.

The waiter task simulates the job of waiter, in which a set of objects needs to be transferred from one location to another distant location. Waiters use a tray to minimize the walked distance. In our simulation, a mobile robot is equipped with a tray-like space that it can use for transferring objects. Discovering the tray in the planning state-space is non-trivial, which makes this task challenging for the planners that do not utilize demonstrations.

B. Evaluation

Baselines. We compare the proposed approach with the following baselines: (i) RRT-connect implemented in HPP [5] without and (ii) with path optimization, (iii) PDDLStreams method [11], and (iv) the proposed approach without and (v) with the proposed path optimization. For PDDLStreams we set the goal requirement to match the final object poses for every object. The object grasp definitions are set to be consistent with the proposed approach.

Metrics. To compare the proposed approach with the baselines, we solve the planning problem on the same machine and use the following metrics: (i) the time to find a solution measured in seconds, (ii) the number of object grasps and releases, *i.e.*, the number of contacts, and (iii) the length of the planned path. The length of the planned path is evaluated as the sum of Euclidean distances computed between the consecutive robot configurations of the planned path.

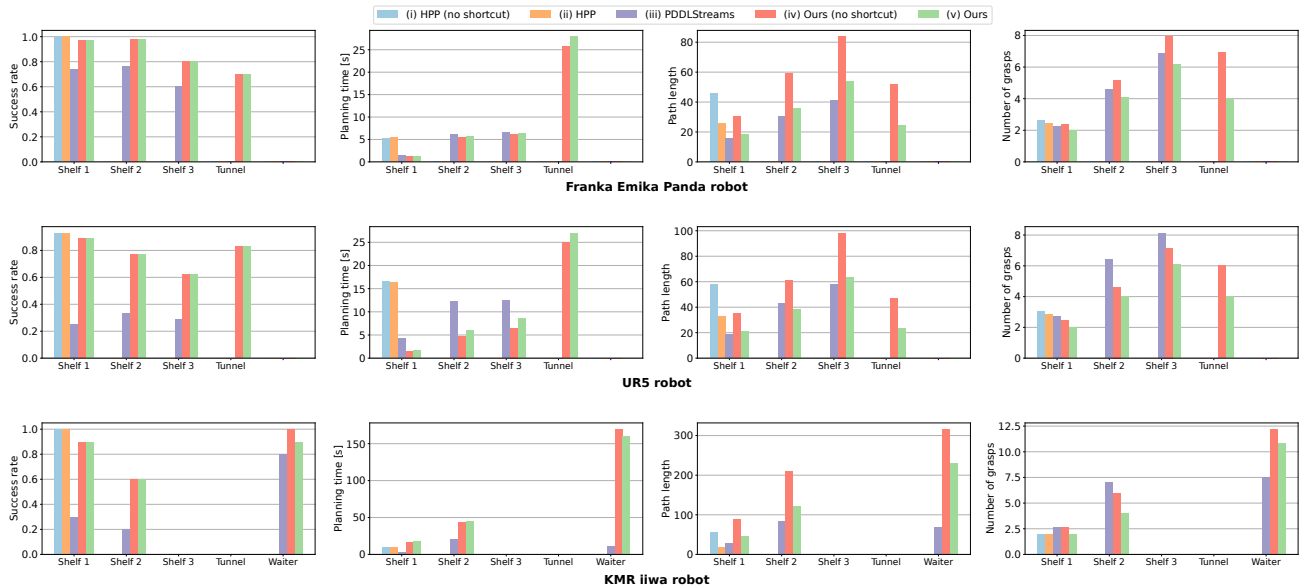


Fig. 4. **Results** reported for different robots (rows) and different metrics (columns). We report (from left): the success rate, the planning time [s], the path length, and the number of grasps. For the success rate, the higher number the better; for the other metrics lower numbers are better. The plots show the comparison of our method with the following baselines: (i) RRT-connect implemented in HPP [5] without and (ii) with the path optimization, (iii) PDDLStreams planner [11], and (iv) our proposed approach without and (v) with the proposed path optimization. Each graph shows results for five different tasks. “Shelf 1” - “Shelf 3” correspond to the Shelf tasks with 1-3 objects. “Tunnel” corresponds to the tunnel task and “Waiter” corresponds to the waiter task. Please note that the Waiter task is only performed with the KMR iiwa mobile robot.

Results are shown in Fig. 4. We set the time limit to 60 s for all tasks except the waiter task where the limit was set to 300 s. Planning was repeated 10 times for KMR robot, and 10 times for each of 10 robot base poses for arm robots. The success rate is reported in the first column of Fig. 4.

For the successful runs, we measure the aforementioned metrics and report the average. The different columns in Fig. 4 show the different reported metrics. The different rows in Fig. 4 correspond to experiments with different robots (from top): (i) 7 DoFs Franka Emika Panda robot, (ii) 6 DoFs UR5 robot, and (iii) 7 DoFs KUKA IIWA arm mounted on 3 DoFs mobile platform (*i.e.* KMR IIWA robot). We sample the robot base poses such that all objects are reachable in both the start and the goal poses. The same robot base poses were used for all the methods. Next, we discuss the results for each task separately.

HPP [5] fails to solve shelf tasks for more than one object in the given time limit. Our method outperforms the PDDLStream planner [11] based on success rate. For the KMR robot, all methods fail to solve the shelf task with three objects and the tunnel task, most likely due to the large base of the robot colliding with the environment. To conclude, the main benefit of guiding the planner by video for the shelf task lies in the success rate, *i.e.* the guidance by video allows us to solve the given task more reliably.

For the tunnel task, the state-of-the-art planners fail to solve the task due to the difficulty of discovering the narrow passage by the sampling-based planning methods. The proposed approach takes advantage of the available demonstration to discover the narrow passage and manages to plan a path for the non-mobile robots.

For the waiter task, we achieve a higher success rate than

the PDDLStream planner [11], however, it takes more time to arrive at a solution. The PDDLStream planner ignores the tray and moves each object one by one, which results in a lower number of grasps.

Limitations and assumptions. Our approach has several assumptions and limitations. We assume that the demonstration consists of a sequence of pick-and-place motions (without pushing or sliding), where we manipulate one object at a time. Start and goal object poses should be reachable by the robot and suitable for grasp. For example, objects should not be too close to each other, otherwise the planning will fail as the gripper is not able to grasp the object. For the waiter task, we assume that objects are manipulated only when the moving platform is stable. For accurate processing of the video demonstration, objects and hands should be well visible at the contact events. Significant occlusions can still pose a challenge for the current approach.

V. CONCLUSIONS

We present an approach for guiding task and motion planning by a video demonstration. We show that contact states and 6D object poses extracted from the video help building the tree in the configuration space of the planner much more efficiently. Our evaluation on our newly built task & motion planning benchmark shows that the video demonstration improves the success rate of planning and is especially effective in environments with narrow passages, which present a major challenge for current sampling based planners. This work makes a step towards large-scale learning of task & motion plans from instructional videos.

REFERENCES

- [1] K. M. Lynch and F. C. Park, *Modern robotics*. Cambridge University Press, 2017.
- [2] L. E. Kavraki, P. Švestka, J.-C. Latombe, and M. Overmars, “Probabilistic roadmaps for fast path planning in high dimensional configuration spaces,” *IEEE Transactions on Robotics and Automation*, vol. 12, no. 4, pp. 566–580, 1996.
- [3] J. Kuffner and S. LaValle, “Rrt-connect: An efficient approach to single-query path planning,” in *Proceedings of 2000 ICRA. Millennium Conference. IEEE International Conference on Robotics and Automation. Symposia Proceedings (Cat. No.00CH37065)*, vol. 2, 2000, pp. 995–1001 vol.2.
- [4] J. Kuffner, K. Nishiwaki, S. Kagami, M. Inaba, and H. Inoue, “Motion planning for humanoid robots,” in *Robotics research. The eleventh international symposium*. Springer, 2005, pp. 365–374.
- [5] F. Lamiroux and J. Mirabel, “Prehensile manipulation planning: Modeling, algorithms and implementation,” *IEEE Transactions on Robotics*, vol. 38, no. 4, pp. 2370–2388, 2022.
- [6] Guided TAMP Project Webpage. [Online]. Available: <https://data.ciirc.cvut.cz/public/projects/2023VideoGuidedTAMP>
- [7] S. M. LaValle *et al.*, “Rapidly-exploring random trees: A new tool for path planning,” *The annual research report*, 1998.
- [8] M. Strandberg, “Augmenting rrt-planners with local trees,” in *IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA '04. 2004*, vol. 4, 2004, pp. 3258–3262 Vol.4.
- [9] J. Janoš, V. Vonásek, and R. Pěnička, “Multi-goal path planning using multiple random trees,” *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 4201–4208, 2021.
- [10] D. Devaurs, T. Siméon, and J. Cortés, “A multi-tree extension of the transition-based rrt: Application to ordering-and-pathfinding problems in continuous cost spaces,” in *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2014, pp. 2991–2996.
- [11] C. R. Garrett, T. Lozano-Pérez, and L. P. Kaelbling, “Pddlstream: Integrating symbolic planners and blackbox samplers via optimistic adaptive planning,” *Proceedings of the International Conference on Automated Planning and Scheduling*, vol. 30, no. 1, pp. 440–448, Jun. 2020. [Online]. Available: <https://ojs.aaai.org/index.php/ICAPS/article/view/6739>
- [12] A. Krontiris and K. E. Bekris, “Dealing with difficult instances of object rearrangement,” in *Robotics: Science and Systems*, vol. 1123, 2015.
- [13] H. Shuai, N. Stiffler, A. Krontiris, K. E. Bekris, and J. Yu, “High-quality tabletop rearrangement with overhead grasps: Hardness results and fast methods,” in *Robotics: Science and Systems (RSS)*, Cambridge, MA, 07/2017 2017. [Online]. Available: <https://arxiv.org/pdf/1705.09180.pdf>
- [14] Y. Labbé, S. Zagoruyko, I. Kalevtykh, I. Laptev, J. Carpentier, M. Aubry, and J. Sivic, “Monte-carlo tree search for efficient visually guided rearrangement planning,” *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3715–3722, 2020.
- [15] G. Ye and R. Alterovitz, “Guided motion planning,” in *Robotics research*. Springer, 2017, pp. 291–307.
- [16] J. Aleotti and S. Caselli, “Part-based robot grasp planning from human demonstration,” in *2011 IEEE International Conference on Robotics and Automation*, 2011, pp. 4554–4560.
- [17] A. Nguyen, D. Kanoulas, L. Muratore, D. G. Caldwell, and N. G. Tsagarakis, “Translating videos to commands for robotic manipulation with deep recurrent neural networks,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 3782–3788.
- [18] D. Shan, J. Geng, M. Shu, and D. F. Fouhey, “Understanding human hands in contact at internet scale,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [19] Y. Labbé, J. Carpentier, M. Aubry, and J. Sivic, “Cosypose: Consistent multi-view multi-object 6d pose estimation,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Cham: Springer International Publishing, 2020, pp. 574–591.
- [20] Z. Cao, I. Radosavovic, A. Kanazawa, and J. Malik, “Reconstructing hand-object interactions in the wild,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 12 417–12 426.
- [21] Y. Hasson, G. Varol, D. Tzionas, I. Kalevtykh, M. J. Black, I. Laptev, and C. Schmid, “Learning joint reconstruction of hands and manipulated objects,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [22] S. Levine, C. Finn, T. Darrell, and P. Abbeel, “End-to-end training of deep visuomotor policies,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1334–1373, 2016.
- [23] M. P. Deisenroth, C. E. Rasmussen, and D. Fox, “Learning to control a low-cost manipulator using data-efficient reinforcement learning,” *Robotics: Science and Systems VII*, vol. 7, pp. 57–64, 2011.
- [24] P. Ennen, P. Bresenitz, R. Vossen, and F. Hees, “Learning robust manipulation skills with guided policy search via generative motor reflexes,” in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 7851–7857.
- [25] J. Nocedal and S. J. Wright, *Numerical Optimization*, 2nd ed. New York, NY, USA: Springer, 2006.
- [26] G. Sánchez and J.-C. Latombe, “On delaying collision checking in prm planning: Application to multi-robot coordination,” *The International Journal of Robotics Research*, vol. 21, no. 1, pp. 5–26, 2002. [Online]. Available: <https://doi.org/10.1177/027836402320556458>
- [27] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, “Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes,” *Robotics: Science and Systems (RSS)*, 2018.