

NOD-TAMP: Multi-Step Manipulation Planning with Neural Object Descriptors

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Developing intelligent robots for complex manipulation tasks in house-
2 hold and factory settings remains challenging due to long-horizon tasks, contact-
3 rich manipulation, and the need to generalize across a wide variety of object
4 shapes and scene layouts. While Task and Motion Planning (TAMP) offers a
5 promising solution, its assumptions such as kinodynamic models limit applicabil-
6 ity in novel contexts. Neural object descriptors (NODs) have shown promise in
7 object and scene generalization but face limitations in addressing broader tasks.
8 Our proposed TAMP-based framework, NOD-TAMP, extracts short manipulation
9 trajectories from a handful of human demonstrations, adapt these trajectories using
10 NOD features, and compose them to solve broad long-horizon tasks. Validated in
11 a simulation environment, NOD-TAMP effectively tackles varied challenges and
12 outperforms existing methods, establishing a cohesive framework for manipula-
13 tion planning. For videos and other supplemental material, see the project website:
14 <https://sites.google.com/view/nod-tamp/>.

15 **Keywords:** Task and Motion Planning; Learning from Demonstration; Neural
16 Object Representations

17 1 Introduction

18 Developing intelligent robots that can automate complex manipulation tasks in households or fac-
19 tories has been a longstanding goal for robotics and AI. Among the multitudes of challenges, three
20 key factors stand out. We illustrate these challenges in a simulated tabletop cleaning task in Fig. 1.
21 First, these tasks are often *long-horizon* and full of sequential dependencies. For example, the robot
22 must reason about the best pose to grasp a mug in order to stow it in a tight cabinet, among other
23 steps to organize the entire table. Second, the *contact-rich* manipulation steps, such as the process
24 of stowing the mug, can make model-based planning intractable [1]. Finally, to be effective in broad
25 environments, the robot must handle a wide *variation of object shapes and scene layouts*.

26 Task and Motion Planning (TAMP) [2, 3] has been the de facto solution for such problems be-
27 cause it can effectively resolve sequential dependencies through hybrid symbolic-continuous plan-
28 ning. However, TAMP systems typically require accurate, special-purpose perception and hand-
29 engineered manipulation skills. This makes it difficult to apply these methods to previously unseen
30 objects and tasks with complex dynamics such as contact-rich manipulation. Recent works have
31 proposed to learn manipulation skills from demonstration [4, 5] to partially relax these constraints.
32 However, their generalization ability remains bounded by the training data, which is costly and
33 difficult to collect at scale [6].

34 At the same time, neural representation models trained on broad data have shown remarkable poten-
35 tial in enabling generalizable manipulation systems [7–10]. In particular, neural object descriptors
36 (NODs) [8, 11, 12] emerged as a powerful tool to extract dense, part-level features that generalize
37 across object instances. Simeonov and Du *et al.* [8] showed that Neural Descriptor Fields (NDF),

38 a type of NOD that encodes SE(3) poses relative to a given object, can adapt key-frame actions
 39 (e.g. grasp poses) for one object instance to others in the same object category (e.g. mugs), thereby
 40 achieving category-level generalization. However, despite this progress, existing NOD-based meth-
 41 ods [8, 13, 14] are limited to adapting individual key-frame poses instead of solving a long-horizon
 42 task, and they struggle with contact-rich manipulation because they still rely on conventional mo-
 43 tion planners to generate the approaching trajectories. Leveraging NODs to solve long-horizon tasks
 44 requires addressing a number of fundamental limitations, including planning long action sequences
 45 and generating trajectories for contact-rich manipulation.

46 To address these limitations, we intro-
 47 duce NOD-TAMP, a TAMP-based
 48 framework that extracts adaptable
 49 skills from a handful of human
 50 demonstrations using NOD features
 51 and compose the skills to solve di-
 52 verse long-horizon tasks. Central
 53 to NOD-TAMP is a skill reasoning
 54 module that composes short-horizon
 55 skills to solve novel long-horizon

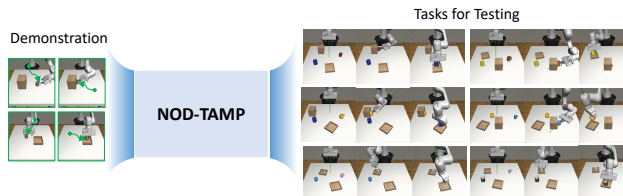


Figure 1: **Overview.** A subset of the diverse long-horizon tasks NOD-TAMP can solve with just a handful of demonstrations.

56 goals never seen in demonstration. To synthesize fine-grained manipulation trajectories and adapt
 57 to new objects, we propose a NOD-based trajectory adaptation module that can consistently adapt a
 58 recorded skill trajectory according to the observed objects. Finally, NOD-TAMP can flexibly switch
 59 between adapting recorded trajectories and using existing kinematic motion planning ability of a
 60 TAMP system to generalize to drastically different scene layout.

61 We evaluate NOD-TAMP with simulated multi-step manipulation tasks that evaluate different fac-
 62 tors of generalization across long-horizon and contact-rich tasks, including object shapes, number
 63 of objects, scene layout, task length, and task objectives. We empirically demonstrate that NOD-
 64 TAMP can consistently solve the evaluation tasks, despite using just a small set of short-horizon
 65 demonstrations. NOD-TAMP also outperforms other methods [8, 15] that share a subset of its traits,
 66 highlighting the value of building a cohesive manipulation planning system.

67 2 Related Work

68 **TAMP.** Task and Motion Planning (TAMP) is a powerful paradigm for addressing long-horizon
 69 manipulation challenges by decomposing a complex planning problem into a series of simpler sub-
 70 problems [2, 3, 16–18]. Nonetheless, TAMP techniques presuppose knowledge of the object models
 71 and the underlying kinodynamic systems. Such presuppositions can be limiting, particularly for
 72 real-world domains with diverse objects and steps that involve complex physical processes such as
 73 contact-rich manipulation.

74 **Learning for TAMP.** Recent works have set to address such limitations by replacing hand-crafted
 75 components in a TAMP system with learned ones. Examples include environment models [19–
 76 22], skill operator models [4, 23], skill samplers [24, 25], and learning to imitate actions from
 77 TAMP supervisors [26–28]. However, these learned components are often limited to the tasks and
 78 environments that they are trained on. Two notable exceptions are MOM [29] and GenTP [30], but
 79 both methods plan with predefined manipulation skills and assume the skills are robust to variations
 80 in tasks and environments. In contrast, our work directly tackles the generalization challenge at
 81 the level of motion generation. Closely related to our work are methods that learn manipulation
 82 skills for TAMP systems [4, 31, 32]. However, the resulting systems remain bottlenecked by the
 83 generalizability of the skills, which are trained using conventional Reinforcement Learning [31] or
 84 Behavior Cloning [4, 32] algorithms. Instead, our work develops skills with object category-level
 85 generalization ability and integrates such skills with the existing planning ability of a TAMP system.

86 **Learning from Human Demonstrations.** Modern deep imitation learning techniques have shown
 87 remarkable performance in solving real-world manipulation tasks [6, 33–37]. However, the promi-

88 nent data-centric view of imitation learning [6, 37, 38], i.e., scaling up robot learning via brute-
 89 force data collection, remains limited by the sample efficiency of the existing learning algorithms
 90 and the challenges in collecting demonstrations for long-horizon tasks in diverse settings. Other
 91 recent works have proposed to replay a small set of human demos in new situations to facilitate
 92 sample-efficient generalization [15, 39–44], but replay without adaptation can fail for novel object
 93 instances. Some other works leverage pretrained object representations to dramatically improve the
 94 generalization of policies given a handful of demonstrations [8, 10, 14]. However, these methods
 95 are limited to adapting a short skill [10] or a single manipulation action [8]. Our work develops
 96 a long-horizon planning framework that seamlessly integrates skills augmented with latent object
 97 representations into a classical TAMP framework.

98 3 Problem Formulation

99 We seek to apply NDFs [45], a type of Neural Object Descriptor, to robot manipulation. First,
 100 we review background on NDF for category-consistent frame estimation (Section 3.1). Then, we
 101 describe our problem setting (Section 3.2).

102 3.1 Neural Descriptor Fields (NDF)

103 NDF [8] were first proposed for category-invariant modeling of object rigid transformations. A NDF
 104 model $F(T | P)$ takes in an object shape that is represented as a point cloud $P \in \mathbb{R}^{N \times 3}$ and a set
 105 of query positions $\{R \cdot x_i + t \mid x_i \in X\}$. Let $T = [R \mid t] \in \text{SE}(3)$ be a rigid transformation
 106 and $X \in \mathbb{R}^{M \times 3}$ be a vector of query points. The NDF outputs features corresponding to the query
 107 points:

$$z \leftarrow F(T | P) \equiv \bigoplus_{x_i \in X} f(R \cdot x_i + t | P).$$

108 A key advantage of NDF features is that they are only related to the queries in the object’s local
 109 frame, therefore a rigid transformation $R \in \text{SE}(3)$ applied to both the shape and the query pose has
 110 no affect on the feature:

$$F(R \cdot T | R \cdot P) = F(T | P).$$

111 NDFs are typically used to solve an inverse problem: recover a transformation T that corresponds
 112 to a query feature z . This can be framed as the following optimization problem and solved by
 113 1) randomly initializing the transform, 2) backpropgating to compute the error gradient, and 3)
 114 iteratively moving along the gradient to minimize the error:

$$\text{NDF-OPTIMIZE}(F, P, z) \equiv \underset{T}{\text{argmin}} \|z - F(T | P)\|.$$

115 3.2 Problem Setting

116 We address controlling a robot to perform multi-stage manipulation. The robot is tasked with manip-
 117 ulating one or more movable objects $o \in O$ in the world to achieve a goal. The robot observes RGB-
 118 D images, which it can process into segmented point clouds for each object $\mathcal{P} = \{o : P_o \mid o \in P\}$.
 119 Although we assume that we can detect the category of each object, critically, we do not assume
 120 instance detection or have a geometric model of any object (*e.g.* mesh).

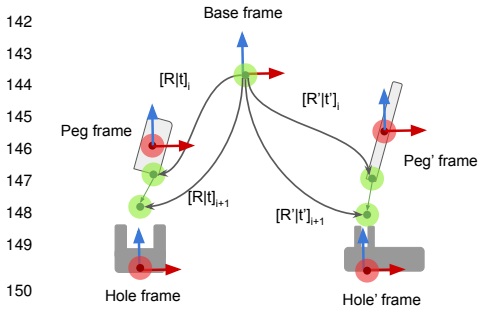
121 The state of the world is described by robot configuration $q \in \mathbb{R}^d$ and frame state $S = \{o : T_w^o \mid$
 122 $o \in O\}$, where $T_{f_2}^{f_1} \in \text{SE}(3)$ represents a rigid transformation from frame f_1 to f_2 and w is the
 123 world frame. Let $S_{f_2}^{f_1}$ be shorthand for the rigid transformation from frame f_1 to f_2 in state S . The
 124 robot takes actions $a = T_a^e$ that correspond to target end-effector poses, where e is the end-effector
 125 frame. These task-space set points are converted to joint-space commands using Operational-Space
 126 Control (OSC) [46].

127 We are interested in re-purposing a set of demonstrations for a ensemble of tasks into manipulation
 128 policy that is effective in scenes with new objects and varying initial poses. To that end, we assume

129 a dataset of demonstrations that is collected by human teleoperation or some other process. Let a
 130 demonstration trajectory $\tau = [\langle S_1, a_1 \rangle, \dots, \langle S_h, a_h \rangle]$ be a sequence of state-action pairs $\langle S, a \rangle$. At
 131 training time, we assume that we can observe or estimate frame states; however, this assumption
 132 does not hold at test time.

133 3.3 Skill Demonstrations

134 In order to deploy demonstrations new 1) objects and 2) tasks, we need to understand more about
 135 the context behind each action involving which objects are interacting or are about to interact as
 136 well as qualitatively how they are interacting. In this work, we assume that each state-action pair
 137 $\langle S, a \rangle$ can be labeled with an interaction type l , a source movable object o , and target coordinate
 138 frame f , forming a tuple $\langle l, o, f, S, a \rangle$. Our technique will be to characterize the motion of o relative
 139 to f using NDFs and then adapt it to new circumstances. Figure 2 demonstrates a pair of insertion
 140 interactions involving pegs and holes that vary in geometry. Here, the movable object o is the peg
 141 and the target frame f is the hole.



152 **Figure 2: Trajectory Adaptation.** Il-
 153 lustration of generating motion trajec-
 154 tory for test scenario based on the ref-
 155 erence demonstration.

156
 157 *feature trajectory* $\mathcal{Z} = [z_1, \dots, z_k]$. And when applied to a segmented demonstration, this produces
 158 a *feature demonstrations* $d = [\langle l_1, o_1, f_1, \mathcal{Z}_1 \rangle, \dots, \langle l_h, o_h, f_h, \mathcal{Z}_h \rangle]$. Thus, we accumulate a *skill*
 159 *dataset* of feature demonstrations $\mathcal{D} = \{d_1, \dots, d_n\}$, which can includes demonstrations that span
 160 multiple tasks.

161 3.4 Peg-in-Hole Example

162 Continuing the ‘‘Peg-in-Hole’’ running example in figure 2, we discuss relevant skills and plans. It
 163 requires two types of interactions: 1) moving the end-effector to grasp an object and 2) inserting a
 164 grasped object into another entity. These interactions can be formalized by the following planning
 165 operators, where $l \in \{\text{GRASP}, \text{ATTACH}\}$:

- 166 1. $\text{GRASP}(o, e; \mathcal{Z})$: move to grasp object o with end-effector frame e using feature trajectory
 167 \mathcal{Z} .
- 168 2. $\text{INSERT}(o, f; \mathcal{Z})$: while grasping object o , move to insert o relative to frame f using feature
 169 trajectory \mathcal{Z} .

170 A plan that directly adapts the demonstrations has the following form:

$$\pi = [\text{GRASP}(\text{peg}, \text{ee}; \mathcal{Z}_1), \text{INSERT}(\text{peg}, \text{hole}; \mathcal{Z}_2)].$$

171 In between the contact-involved component of these interactions are contact-adverse robot move-
 172 ment. We can plan motions between the end of the last component and the start of the next one in
 173 order to more robustly and efficiently move between segments. These segments can also be repre-
 174 sented by planning operators [47]:

- 175 1. TRANSIT(τ): while not grasping any object, move the robot along trajectory τ .
 176 2. TRANSFER($o; \tau$): while grasping object o , move the robot along trajectory τ .

177 A plan that includes motion operators has the following form:

$$\pi = [\text{TRANSIT}(\tau_1), \text{GRASP}(\text{peg}, \text{ee}; \mathcal{Z}_1), \\ \text{TRANSFER}(\text{peg}; \tau_2), \text{ATTACH}(\text{peg}, \text{hole}; \mathcal{Z}_2)]$$

178 4 NOD-TAMP

179 We present a set of algorithms for
 180 adapting a dataset of demonstrations
 181 to new problems. First, we show
 182 how a single demonstration can be
 183 adapted to a new environment using
 184 NDFs (Section 4.1). Then, we pro-
 185 pose a planning algorithm that’s able
 186 to combine skill segments from multiple
 187 demonstrations to maximize effec-
 188 tiveness (Section 4.2). Finally, we
 189 use motion planning to connect each segment in order to efficiently and robustly generalize to new
 190 workspaces (Section 4.3).

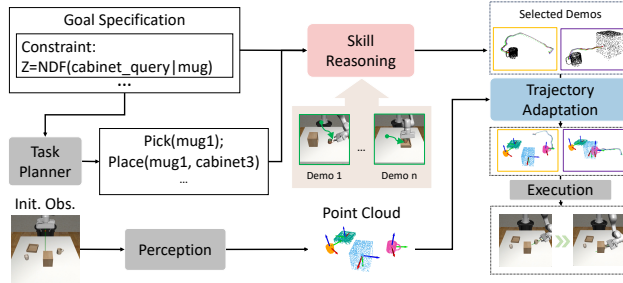
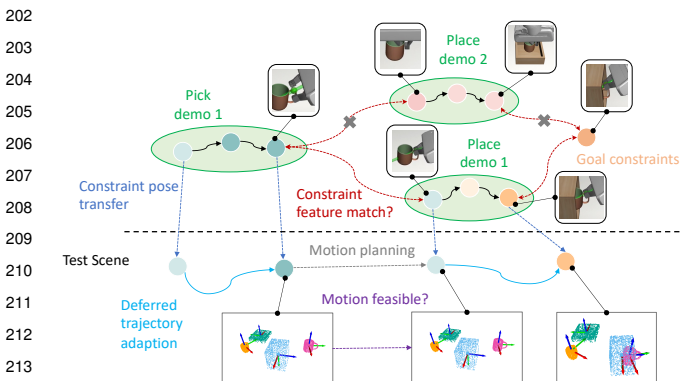


Figure 3: **System.** Overview of the NOD-TAMP system.

191 4.1 Trajectory Adaptation

192 The first algorithm we introduce directly adapts a demonstration d to the current task. Algorithm 1
 193 displays the trajectory adaptation pseudocode. It iterates over each labeled feature trajectory in
 194 demonstration d and then over each timestep in the trajectories. For each timestep, it computes
 195 the target transformation T_z that corresponds to NDF feature z . Depending on whether object o is
 196 attached to another frame f' in the scene graph G , it composes the transform into a target world pose
 197 T_w^f for manipulation frame f . It then converts this into a target end-effector pose T_w^f , which will
 198 serve as a setpoint for a downstream controller or planner. After iterating over a labeled trajectory,
 199 the scene graph G is updated with the new state of the manipulated object o . This also models that
 200 the point cloud P_o for object o is now attached to and thus moving with frame f .

201 4.2 Skill Planning



214 Figure 4: **Skill Planning and Trajectory Generation.** We
 215 illustrate how skill planning and trajectory adaptation col-
 216 laborated to generate mug insertion trajectories in this ex-
 217 ample.
 218

It can be limiting to be only able to reuse whole demonstrations. First, it is inflexible in problems where the, *plan skeleton*, or the sequence of skills changes. Second, segments of multiple demonstrations in conjunction might better address a new problem. For example, we may just transfer the picking part of the pick-and-place trajectory that work with a mug to interact with a new object, e.g., a bowl.

Algorithm 2 displays the pseudocode for the NOD-TAMP planner. It takes in a plan skeleton $\hat{\pi}$ that defines a sequence of skill types to consider. It first compiles a list of skills in dataset

219 \mathcal{D} relevant to π . Then, it iterates through all possible plans using these skills. Each candidate π is
 220 scored based on the compatibility of subsequent actions using NDF features, the score c is computed
 221 as $c = \sum_{i=1}^{|\pi|-1} \|z_{i+1} - z_i\|$, where for specific i , the NDF feature z_i and z_{i+1} are extracted as:

$$z_i, z_{i+1} \leftarrow F(T_i | P), F(T_{i+1} | P)$$

222 Here T_i and T_{i+1} are the query pose determined at the last time step and first time step of the
 223 trajectories from skill i and skill $i + 1$ respectively, and P is the point cloud of the target object for
 224 skill i and skill $i + 1$ captured during demonstration.

225 After we obtain all scores for all skill combinations, the plan with the lowest plan-wide NDF feature
 226 distance is returned. For simplicity, we present this as a Cartesian product over relevant skills, but
 227 this can be done more efficiently by performing a Uniform Cost Search in plan space, where the
 228 NDF feature distance serves as the cost function. The visualization of the feature matching process
 229 for a mug placing example is presented in Fig. 5, showing that the grasping on the mug rim would be
 230 compatible with placing the mug uprightly in a bin, and grasp on the handle would be compatible
 231 with inserting the mug into a cabinet.

232 4.3 Transit & Transfer Motion

233 Adapting demonstrated skills is particularly effective at generating behavior that involves contact.
 234 However, demonstrations typically contain long segments without contact (outside of holding an object).
 235 Because these components do not modify the world, it is often not productive to replicate them.
 236 Because these components do not modify the world, it is often not productive to replicate them.
 237 Thus, we temporally trim skill demonstrations to focus on the data points that involve contact.
 238 In our implementation, we simply select the 50 steps that are most close to the time point when the
 239 contact happening.

244 After trimming, and in many cases before trimming, two adjacent skills might be quite far away in task space.
 245 While we could simply interpolate between them, this is not generally safe because the straight-line path may cause the robot to unexpectedly collide.
 246 To address this, we use motion planning to optimize for safe and efficient motion that reaches the start of next the skill.
 247 Motion planning generally requires some characterization of the collision volume of the obstacles to avoid. Because we do not assume access to object models,
 248 we use the segmented point clouds \mathcal{P} the collision representation, where each point in the cloud is a sphere with radius $r > 0$.

252 Algorithm 3 displays the pseudocode for the full NOD-TAMP policy with motion planning. It displayed in a manner with online motion planning and execution, but full plans can also be computed offline.
 253 The policy first computes a skill plan π and then adapts each labeled trajectory in sequence. Between trajectories, it uses PLAN-MOTION to plan a trajectory from the current robot configuration
 254 q to a configuration that reaches the first end-effector pose T_w^e . Specifically, it samples goal configurations q' using inverse kinematics and then invokes a sampling-based motion planner between q
 255 and q' . Once the end-effector arrives at T_w^e , for each pose yielded by the skill, the policy deploys OSC to track the target poses.

260 A detailed example of how the system work for placing the mug into the cabinet is presented in Fig. 4. The skill planning process first identifies the useful trajectories of the pick skill and place skill through matching the constraint feature.
 261 A motion planner then generates the transition motion between skills. Once the motion transition among every connections between the consecutive skills
 262 are determined, we transfer the trajectories to the test scene through trajectory adaptation.

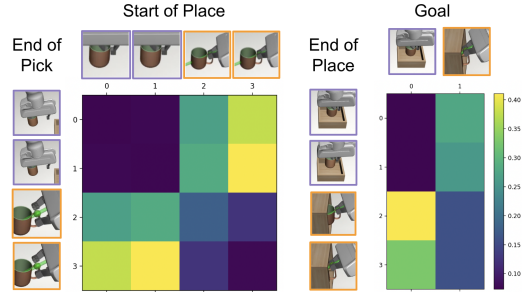


Figure 5: **Feature Matching.** We show the feature matching distance for different skill combinations when placing the mug.

265 5 Experiments

266 We aim to validate NOD-TAMP and how its components contribute to solve long-horizon tasks,
267 perform contact-rich manipulation, and generalize to new object shapes.

268 5.1 Experimental Setup

269 Tasks.

270 We introduce four simulated [48] table organization
271 tasks (Fig. 6 bottom), which vary in diffi-
272 culty and generalization challenges. We use 10
273 mug models of varying shapes and dimensions
274 from the ShapeNet dataset [49]. Demonstra-
275 tions are provided for one mug and the testing
276 environment randomly samples among the re-
277 maining 9 mugs. We provide all methods with
278 only **two** picking and placing skill trajectories
279 on **one** mug instance (Fig. 6 top). The two tra-
280 jectories vary by grasping pose (side vs. top)
281 and placing pose (side vs. top).

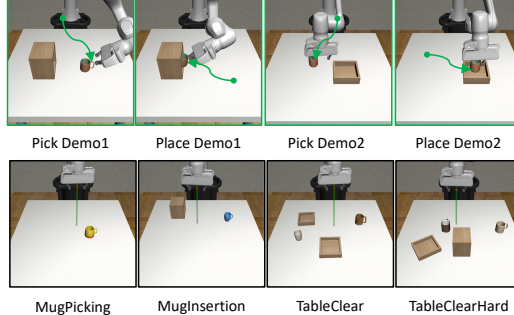


Figure 6: **Demonstration Skills and Tasks.** Il-
lustration of the demonstration skills and the task
setups.

282 **MugPicking** - Pick up different mugs with varying shapes and initial poses. The task tests the ability
283 to adapt the trajectory based on object shapes and poses.

284 **MugInsertion** - Insert mugs of varying shape into a tight cabinet. Both the mug and the cabinet
285 are randomly placed on the table. This task requires fine-grained manipulation for the insertion and
286 adaptation to different initial configurations.

287 **TableClear** - This long-horizon task requires the robot to place two mugs into two bins, which aims
288 to test the ability to achieve long-term goals by reusing the skills.

289 **TableClearHard** - This long-horizon task requires the robot to stow one mug into a cabinet with
290 side opening and place another mug into a bin. The robot must select a feasible chain of skills
291 (side-picking to side-stowing) to transport each mug.

292 **Baselines.** We compare NOD-TAMP with ablation baselines and adapt state-of-the-art methods to
293 facilitate fair comparison. All NDF-based methods share the same NDF model pretrained on mug
294 shapes provided by the original implementation [8].

295 **NDF⁺** [8]: We augment the original NDF manipulation planner, which only generates key-frame
296 manipulation poses, with our task planner and the skill reasoning module, which provides the base-
297 line with the target object and the query pose at different stages. This baseline also uses a motion
298 planner to transition between key-frame poses.

299 **MimicGen⁺** [15]: MimicGen directly pieces together contact-rich segments from human demon-
300 strations and linearly interpolates the intermediate steps. The robot control poses in the contact-rich
301 segments are transformed to the relevant object frame and then sent to the controller without further
302 adaptation. To ensure fair comparison, we augment MimicGen with a motion planner for collision
303 avoidance.

304 **Ours w/o Skill Reasoning (Ours-SR):** This ablation removes the skill reasoning module. For each
305 skill, we randomly choose a reference trajectory from the collected demonstrations belonging to this
306 skill, and we bridge the intermediate transition with the motion planner. This baseline validates the
307 importance of skill reasoning for generalizing across tasks.

308 **Ours w/o Motion Planning (Ours-MP):** This ablation removes the motion planning component
309 and uses linear trajectory interpolation to transition between the adapted skill trajectories generated
310 by the optimization. We set up this baseline to validate the benefit of leveraging motion planning, a
311 capability present in TAMP systems.

312 **Naive Skill Chaining (NSC):** This baseline ablates both the skill reasoning and the motion planning
 313 component: it randomly selects a reference trajectory for each skill, adapts the skill with NDF, and
 314 uses linear trajectory interpolation to transition between the selected trajectories.

315 5.2 Main Results

316 Table 1: Evaluation results (success rate) of all
 317 methods.

Tasks	NDF ⁺	MimicGen ⁺	NSC	Ours (-MP)	Ours (-SR)	Ours
MugPicking	80	70	85	80	85	85
MugInsertion	75	55	80	85	80	90
TableClear	60	75	80	75	85	85
TableClearHard	40	55	15	50	10	80

318
 319
 320
 321
 322

323 tion, skill planning, and motion planning components of NOD-TAMP.

324 **NOD-TAMP exhibits strong performance across long-horizon tasks and is able to reuse skills
 325 in new contexts.** The TableClear task requires methods to re-use the existing **two** pick-and-place
 326 human demonstrations, which only consisted of single mug and bin interactions, to stow two mugs
 327 into two bins. NOD-TAMP achieves strong performance (85%) and outperforms MimicGen⁺ by
 328 15% and NDF⁺ by 25% on this task, showcasing a superior ability to re-purposing short-horizon
 329 skill demonstrations for long-horizon manipulation.

330 **NOD-TAMP exhibits strong generalization capability across goals, objects, and scenes in long-
 331 horizon tasks.** The TableClearHard task requires intelligent selection and application of different
 332 demonstration trajectories to achieve different kinds of mug placements. The task also requires deal-
 333 ing with novel mug objects, and novel scene variations. Here, NOD-TAMP achieves 80%, outper-
 334 forming all other baselines by a wide margin (40% better than NDF⁺, 25% better than MimicGen⁺).
 335 We also see the clear benefit of the skill reasoning component to achieve the different goals in this
 336 task – NOD-TAMP outperforms (Ours-SR) by 70% and NSC by 65%. The omission of the skill
 337 reasoning module results in an incompatible composition of skills. For example, the robot may grip
 338 the rim of a mug and attempt to insert it into the cabinet, leading to collisions between the cabinet
 339 and the gripper. Finally, the motion planning component is also valuable in dealing with the scene
 340 variation – NOD-TAMP outperforms (Ours-MP) by 30%. Simply connecting end-effector trajecto-
 341 ries through linear interpolation without considering obstacles often leads to collisions between the
 342 robot or the held object and its surroundings. In particular, we frequently observed such failures
 343 for the Ours (-MP) approach that the gripper became obstructed by the cabinet after completing the
 344 insertion of the first mug.

345 6 Conclusion

346 We introduced NOD-TAMP, a frame-
 347 work for adaptable manipulation
 348 planning using human demonstra-
 349 tions for long-horizon tasks. While
 350 powerful, NOD-TAMP has limita-
 351 tions that inspire future work. First,
 352 because of the expensive NDF-based
 353 pose optimization process, NOD-
 354 TAMP is slow in practice and is far from real-time. We plan to experiment with more efficient
 355 NOD variants and faster optimization procedure. Moreover, an important aspect of TAMP problem
 356 is representing and satisfying constraints. We plan to explore NOD-based constraint representations
 357 and incorporate constraints into the planning objectives. Finally, NOD-TAMP solves the high-level
 358 task plans and low-level trajectories in silos. As a next step, we aim to enable bi-level communica-
 359 tion in planning and extend NOD-TAMP into a full-fledged integrated TAMP system [17].

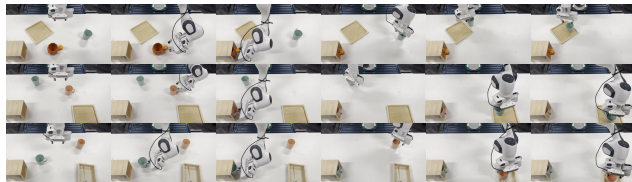


Figure 7: **Real Robot Executions.** Key frames of three task execution processes.

References

- [1] B. Aceituno-Cabezas and A. Rodriguez, “A global quasi-dynamic model for contact-trajectory optimization in manipulation,” 2020.
- [2] L. P. Kaelbling and T. Lozano-Pérez, “Hierarchical task and motion planning in the now,” in *ICRA*, 2011.
- [3] C. R. Garrett *et al.*, “Integrated task and motion planning,” *Annual review of control, robotics, and autonomous systems*, vol. 4, pp. 265–293, 2021.
- [4] T. Silver, R. Chitnis, J. Tenenbaum, L. P. Kaelbling, and T. Lozano-Pérez, “Learning symbolic operators for task and motion planning,” in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE, 2021, pp. 3182–3189.
- [5] S. Cheng and D. Xu, “League: Guided skill learning and abstraction for long-horizon manipulation,” *IEEE Robotics and Automation Letters*, 2023.
- [6] A. Brohan *et al.*, “Rt-1: Robotics transformer for real-world control at scale,” in *arXiv preprint arXiv:2212.06817*, 2022.
- [7] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei, “Voxposer: Composable 3d value maps for robotic manipulation with language models,” *arXiv preprint arXiv:2307.05973*, 2023.
- [8] A. Simeonov *et al.*, “Neural descriptor fields: Se(3)-equivariant object representations for manipulation,” 2022.
- [9] M. Shridhar, L. Manuelli, and D. Fox, “Cliport: What and where pathways for robotic manipulation,” in *Conference on Robot Learning*, PMLR, 2022, pp. 894–906.
- [10] B. Wen, W. Lian, K. Bekris, and S. Schaal, “You only demonstrate once: Category-level manipulation from single visual demonstration,” *arXiv preprint arXiv:2201.12716*, 2022.
- [11] P. R. Florence, L. Manuelli, and R. Tedrake, “Dense object nets: Learning dense visual object descriptors by and for robotic manipulation,” *arXiv preprint arXiv:1806.08756*, 2018.
- [12] P. Sundaresan *et al.*, “Learning rope manipulation policies using dense object descriptors trained on synthetic depth data,” in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2020, pp. 9411–9418.
- [13] A. Simeonov *et al.*, “Se (3)-equivariant relational rearrangement with neural descriptor fields,” in *Conference on Robot Learning*, PMLR, 2023, pp. 835–846.
- [14] E. Chun, Y. Du, A. Simeonov, T. Lozano-Perez, and L. Kaelbling, “Local neural descriptor fields: Locally conditioned object representations for manipulation,” *arXiv preprint arXiv:2302.03573*, 2023.
- [15] A. Mandlekar *et al.*, “Mimicgen: A data generation system for scalable robot learning using human demonstrations,” in *Conference on Robot Learning (CoRL)*, 2023.
- [16] L. P. Kaelbling and T. Lozano-Pérez, “Pre-image backchaining in belief space for mobile manipulation,” in *Robotics Research*, Springer, 2017, pp. 383–400.
- [17] C. R. Garrett, T. Lozano-Pérez, and L. P. Kaelbling, “Pddlstream: Integrating symbolic planners and blackbox samplers via optimistic adaptive planning,” in *Proceedings of the International Conference on Automated Planning and Scheduling*, vol. 30, 2020, pp. 440–448.
- [18] M. A. Toussaint, K. R. Allen, K. A. Smith, and J. B. Tenenbaum, “Differentiable physics and stable modes for tool-use and manipulation planning,” 2018.
- [19] G. Konidaris, L. P. Kaelbling, and T. Lozano-Perez, “From skills to symbols: Learning symbolic representations for abstract high-level planning,” *Journal of Artificial Intelligence Research*, vol. 61, pp. 215–289, 2018.
- [20] Z. Wang, C. R. Garrett, L. P. Kaelbling, and T. Lozano-Pérez, “Learning compositional models of robot skills for task and motion planning,” *The International Journal of Robotics Research*, vol. 40, no. 6-7, pp. 866–894, 2021.
- [21] J. Liang, M. Sharma, A. LaGrassa, S. Vats, S. Saxena, and O. Kroemer, “Search-based task planning with learned skill effect models for lifelong robotic manipulation,” in *2022 International Conference on Robotics and Automation (ICRA)*, IEEE, 2022, pp. 6351–6357.

- 411 [22] A. Simeonov *et al.*, “A long horizon planning framework for manipulating rigid pointcloud
412 objects,” in *Conference on Robot Learning*, PMLR, 2021, pp. 1582–1601.
- 413 [23] H. M. Pasula, L. S. Zettlemoyer, and L. P. Kaelbling, “Learning symbolic models of stochastic
414 domains,” *Journal of Artificial Intelligence Research*, vol. 29, pp. 309–352, 2007.
- 415 [24] R. Chitnis *et al.*, “Guided search for task and motion plans using learned heuristics,” in *ICRA*,
416 IEEE, 2016.
- 417 [25] B. Kim, L. Shimanuki, L. P. Kaelbling, and T. Lozano-Pérez, “Representation, learning, and
418 planning algorithms for geometric task and motion planning,” *IJRR*, vol. 41, no. 2, 2022.
- 419 [26] M. J. McDonald and D. Hadfield-Menell, “Guided imitation of task and motion planning,” in
420 *Conference on Robot Learning*, PMLR, 2022, pp. 630–640.
- 421 [27] J. Gu *et al.*, “Maniskill2: A unified benchmark for generalizable manipulation skills,” *arXiv*
422 *preprint arXiv:2302.04659*, 2023.
- 423 [28] M. Dalal, A. Mandlekar, C. Garrett, A. Handa, R. Salakhutdinov, and D. Fox, “Imitating task
424 and motion planning with visuomotor transformers,” *arXiv preprint arXiv:2305.16309*, 2023.
- 425 [29] A. Curtis, X. Fang, L. P. Kaelbling, T. Lozano-Pérez, and C. R. Garrett, “Long-horizon ma-
426 nipulation of unknown objects via task and motion planning with estimated affordances,” in
427 *2022 International Conference on Robotics and Automation (ICRA)*, IEEE, 2022, pp. 1940–
428 1946.
- 429 [30] C. Wang, D. Xu, and L. Fei-Fei, “Generalizable task planning through representation pre-
430 training,” *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 8299–8306, 2022.
- 431 [31] S. Cheng and D. Xu, “Guided skill learning and abstraction for long-horizon manipulation,”
432 *arXiv preprint arXiv:2210.12631*, 2022.
- 433 [32] A. Mandlekar, C. R. Garrett, D. Xu, and D. Fox, “Human-in-the-loop task and motion plan-
434 ning for imitation learning,” in *7th Annual Conference on Robot Learning*, 2023.
- 435 [33] T. Zhang, Z. McCarthy, O. Jow, D. Lee, K. Goldberg, and P. Abbeel, “Deep imitation
436 learning for complex manipulation tasks from virtual reality teleoperation,” *arXiv preprint*
437 *arXiv:1710.04615*, 2017.
- 438 [34] A. Mandlekar *et al.*, “What matters in learning from offline human demonstrations for robot
439 manipulation,” in *Conference on Robot Learning (CoRL)*, 2021.
- 440 [35] A. Mandlekar, D. Xu, R. Martín-Martín, S. Savarese, and L. Fei-Fei, “Learning to generalize
441 across long-horizon tasks from human demonstrations,” *arXiv preprint arXiv:2003.06085*,
442 2020.
- 443 [36] E. Jang *et al.*, “Bc-z: Zero-shot task generalization with robotic imitation learning,” in *Con-*
444 *ference on Robot Learning*, PMLR, 2022, pp. 991–1002.
- 445 [37] M. Ahn *et al.*, “Do as i can, not as i say: Grounding language in robotic affordances,” *arXiv*
446 *preprint arXiv:2204.01691*, 2022.
- 447 [38] C. Lynch *et al.*, “Learning latent plans from play,” in *Conference on robot learning*, PMLR,
448 2020, pp. 1113–1132.
- 449 [39] N. Di Palo and E. Johns, “Learning multi-stage tasks with one demonstration via self-replay,”
450 in *Conference on Robot Learning*, PMLR, 2022, pp. 1180–1189.
- 451 [40] E. Johns, “Coarse-to-Fine Imitation Learning: Robot Manipulation from a Single Demonstra-
452 tion,” *ICRA*, 2021.
- 453 [41] V. Vosylius and E. Johns, “Where to start? transferring simple skills to complex environ-
454 ments,” *arXiv preprint arXiv:2212.06111*, 2022.
- 455 [42] A. Chenu, O. Serris, O. Sigaud, and N. Perrin-Gilbert, “Leveraging sequentiality in reinforc-
456 e-ment learning from a single demonstration,” *arXiv preprint arXiv:2211.04786*, 2022.
- 457 [43] E. Valassakis, G. Papagiannis, N. Di Palo, and E. Johns, “Demonstrate once, imitate imme-
458 diately (dome): Learning visual servoing for one-shot imitation learning,” in *2022 IEEE/RSJ*
459 *International Conference on Intelligent Robots and Systems (IROS)*, IEEE, 2022, pp. 8614–
460 8621.

- 461 [44] J. Liang, B. Wen, K. Bekris, and A. Boularias, "Learning sensorimotor primitives of sequen-
462 tial manipulation tasks from visual demonstrations," in *2022 International Conference on*
463 *Robotics and Automation (ICRA)*, IEEE, 2022, pp. 8591–8597.
- 464 [45] M. Sundermeyer, A. Mousavian, R. Triebel, and D. Fox, "Contact-graspnet: Efficient 6-dof
465 grasp generation in cluttered scenes," in *2021 IEEE International Conference on Robotics*
466 *and Automation (ICRA)*, IEEE, 2021, pp. 13 438–13 444.
- 467 [46] O. Khatib, "A unified approach for motion and force control of robot manipulators: The
468 operational space formulation," *IEEE Journal on Robotics and Automation*, vol. 3, no. 1,
469 pp. 43–53, 1987.
- 470 [47] T. Siméon, J.-P. Laumond, J. Cortés, and A. Sahbani, "Manipulation planning with proba-
471 bilistic roadmaps," *The International Journal of Robotics Research*, vol. 23, no. 7-8, pp. 729–
472 746, 2004.
- 473 [48] E. Todorov, T. Erez, and Y. Tassa, "MuJoCo: A physics engine for model-based control," Oct.
474 2012, pp. 5026–5033.
- 475 [49] A. X. Chang *et al.*, "Shapenet: An information-rich 3d model repository," *arXiv preprint*
476 *arXiv:1512.03012*, 2015.

477 **7 Appendix**

478 **7.1 Algorithms**

479 We provide the pseudo-code of our proposed algorithms.

Algorithm 1 Trajectory adaptation

Declare: Segmented point clouds \mathcal{P}
Declare: Initial end-effector pose S_w^e
Declare: Demonstration dataset $\mathcal{D} = \{d_1, \dots, d_n\}$
Declare: NDFs \mathcal{F}

- 1: **procedure** ADAPT-TRAJ($\mathcal{P}, S_w^e, d; \mathcal{F}$)
- 2: $\mathcal{S} = \{f : \text{NDF-ESTIMATE}(\mathcal{P}, f) \mid f \in F\} \cup \{e : S_w^e\}$
- 3: $G \leftarrow \{\}$ ▷ Object scene graph
- 4: **for** $\langle l, o, f, \mathcal{Z} \rangle \in d$ **do**
- 5: **for** $z \in \mathcal{Z}$ **do**
- 6: $T_z \leftarrow \text{NDF-OPTIMIZE}(\mathcal{F}[o], \mathcal{P}[o], z)$
- 7: $S_w^e, S_w^f \leftarrow \mathcal{S}[e], \mathcal{S}[f]$ ▷ End-effector & frame
- 8: **if** $o \in G$ **then** ▷ Relative to scene graph
- 9: $\langle f', T_w^{f'} \rangle \leftarrow G[o]$
- 10: $S_w^{f'} \leftarrow \mathcal{S}[f']$
- 11: $T_w^f \leftarrow S_w^{f'} \cdot (T_w^{f'})^{-1} \cdot T_z$
- 12: **else** ▷ Relative to world frame
- 13: $T_w^f \leftarrow T_z$
- 14: $T_w^e \leftarrow S_w^e \cdot (S_w^f)^{-1} \cdot T_w^f$ ▷ End-effector target
- 15: ▷ If $f = e$, this reduces to $T_w^e \leftarrow T_w^f$
- 16: **yield** $\langle l, o, f, T_w^e \rangle$ ▷ Yield target to controller
- 17: $\mathcal{S}[e] \leftarrow T_w^e$
- 18: $G[o] \leftarrow \langle f, S_w^f \rangle$ ▷ Set f as the parent of o

Algorithm 2 NOD-TAMP planner

Declare: Plan skeleton $\hat{\pi} = [\langle l_1, o_1, f_1 \rangle, \dots, \langle l_h, o_h, o_h \rangle]$

- 1: **procedure** PLAN-NDF-SKILLS($\mathcal{P}, \hat{\pi}; \mathcal{D}, \mathcal{F}$)
- 2: $D = []$ ▷ List of demos per skill
- 3: **for** $\langle l_i, o_i, f_i \rangle \in \hat{\pi}$ **do**
- 4: $D \leftarrow D + [\{\langle l, o, f, \mathcal{Z} \rangle \mid d \in \mathcal{D}, \langle l, o, f, \mathcal{Z} \rangle \in d, l=l_i \wedge o=o_i \wedge f=f_i\}]$
- 5: $\pi_* \leftarrow \text{None}; c_* \leftarrow \infty$
- 6: **for** $\pi \in \text{PRODUCT}(D_\pi)$ **do** ▷ All combinations
- 7: $c_\pi \leftarrow 0$ ▷ Feature cost
- 8: **for** $i \in [1, \dots, |\pi|-1]$ **do**
- 9: $\langle l_i, o_i, f_i, \mathcal{Z}_i \rangle \leftarrow \pi[i]$
- 10: $\langle l_{i+1}, o_{i+1}, f_{i+1}, \mathcal{Z}_{i+1} \rangle \leftarrow \pi[i+1]$
- 11: **if** $o_i = o_{i+1}$ **then** ▷ Actions with same object
- 12: $F \leftarrow \mathcal{F}[o]; P \leftarrow \mathcal{P}[o]$
- 13: $T_i \leftarrow \text{NDF-OPTIMIZE}(F, P, \mathcal{Z}_i[-1])$
- 14: $T_{i+1} \leftarrow \text{NDF-OPTIMIZE}(F, P, \mathcal{Z}_{i+1}[0])$
- 15: $z_i, z_{i+1} \leftarrow F(T_i \mid P), F(T_{i+1} \mid P)$
- 16: $c_\pi \leftarrow c_\pi + \|z_{i+1} - z_i\|$
- 17: **if** $c_\pi < c_*$ **then** ▷ Update best plan
- 18: $\pi_* \leftarrow \pi; c_* \leftarrow c_\pi$
- 19: **return** π_*

Algorithm 3 NOD-TAMP policy

```
1: procedure NOD-TAMP-POLICY( $\mathcal{P}, \pi; \mathcal{D}, \mathcal{F}$ )
2:    $\pi \leftarrow$  PLAN-NDF-SKILLS( $\mathcal{P}, \pi; \mathcal{D}, \mathcal{F}$ )
3:   if  $\pi = \text{None}$  then
4:     return False ▷ Skill planning failed
5:    $a \leftarrow \text{None}$  ▷ Current action
6:    $S_w^e \leftarrow$  FORWARD-KIN( $q$ )
7:   for  $\langle o, f, T_w^e \rangle \in$  ADAPT-TRAJ( $\mathcal{P}, S_w^e, \pi; \mathcal{F}$ ) do
8:     if  $\langle o, f \rangle \neq a$  then ▷ Action changed
9:        $a \leftarrow \langle o, f \rangle$  ▷ Update current action
10:       $q \leftarrow$  OBSERVE-CONF()
11:       $\tau \leftarrow$  PLAN-MOTION( $\mathcal{P}, q, T_w^e$ )
12:      if  $\tau = \text{None}$  then
13:        return False ▷ Motion planning failed
14:      EXECUTE-TRAJ( $\tau$ )
15:       $q \leftarrow$  OBSERVE-CONF()
16:      EXECUTE-OSC( $q, T_w^e$ ) ▷ Operational Space
17:   return True ▷ Policy succeeded
```

480 **7.2 More Qualitative Results**

481 With **TWO** picking and placing trajectories on just **ONE** mug, we evaluate our method’s effec-
482 tiveness across diverse tasks with different mug shapes, poses, and goal setups. NOD-TAMP con-
483 sistently achieves a high success rate (80-90%) across all tasks, some executions are visualized in
484 Fig. 8.

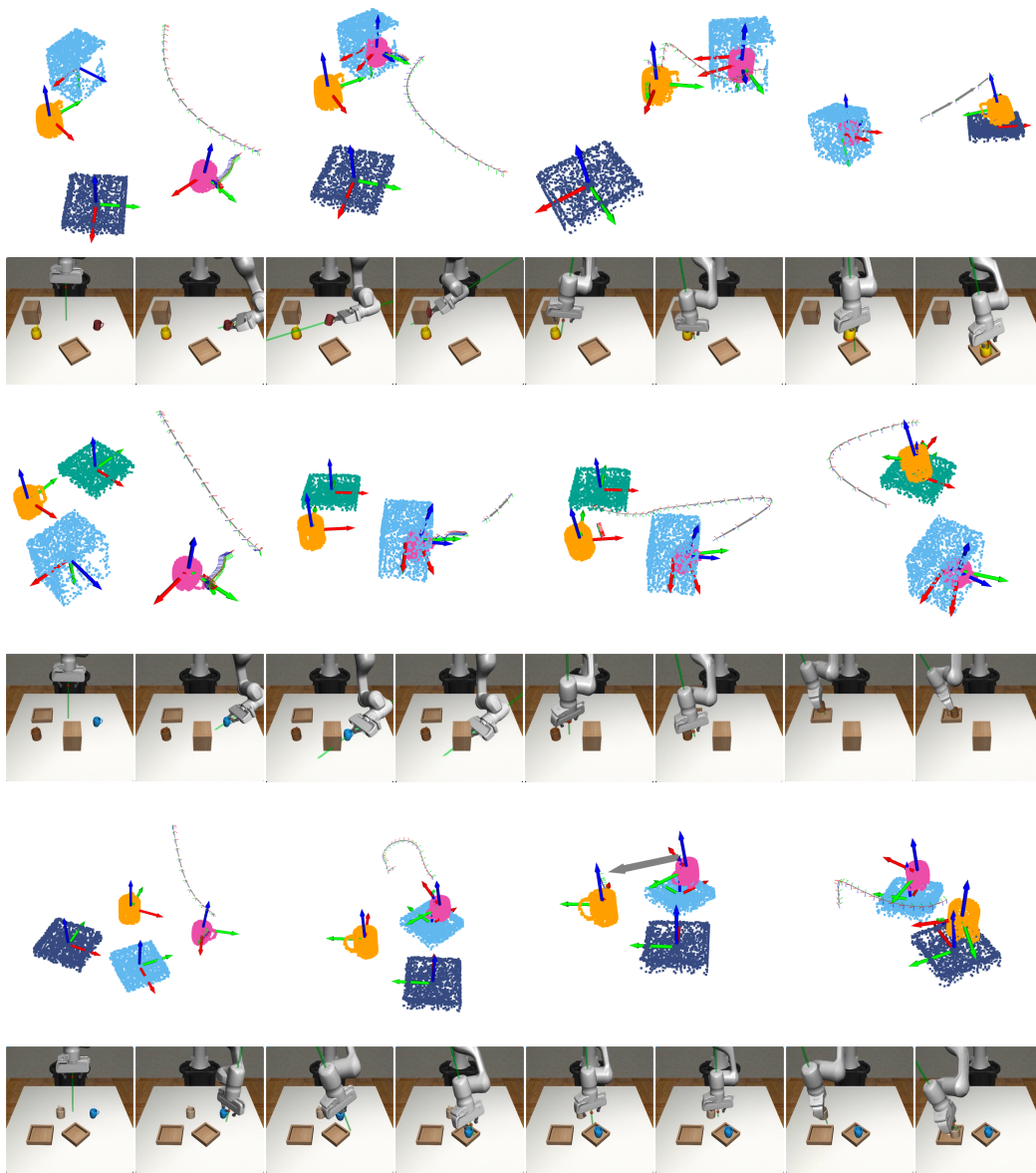


Figure 8: **Qualitative Visualization.** Trajectories generated by our system at different stages, the planned scene represented as point cloud, and snapshots of the execution process.