

ReKep: Spatio-Temporal Reasoning of Relational Keypoint Constraints for Robotic Manipulation

Anonymous Author(s)

1 **Abstract:** Representing robotic manipulation tasks as constraints that associate
2 the robot and the environment is a promising way to encode desired robot be-
3 haviors. However, it remains unclear how to formulate the constraints such that
4 they are 1) versatile to diverse tasks, 2) free of manual labeling, and 3) optimiz-
5 able by off-the-shelf solvers to produce robot actions in real-time. In this work, we
6 introduce **Relational Keypoint Constraints (ReKep)**, a visually-grounded repre-
7 sentation for constraints in robotic manipulation. Specifically, ReKep is expressed
8 as Python functions mapping a set of 3D keypoints in the environment to a numeri-
9 cal cost. We demonstrate that by representing a manipulation task as a sequence
10 of Relational Keypoint Constraints, we can employ a hierarchical optimization
11 procedure to solve for robot actions (represented by a sequence of end-effector
12 poses in $SE(3)$) with a perception-action loop at a real-time frequency. Further-
13 more, in order to circumvent the need for manual specification of ReKep for each
14 new task, we devise an automated procedure that leverages large vision models
15 and vision-language models to produce ReKep from free-form language instruc-
16 tions and RGB-D observations. We present system implementations on a wheeled
17 single-arm platform and a stationary dual-arm platform that can perform a large
18 variety of manipulation tasks, featuring multi-stage, in-the-wild, bimanual, and
19 reactive behaviors, all without task-specific data or environment models. Videos
20 and code can be found at rekep-cornell.github.io.

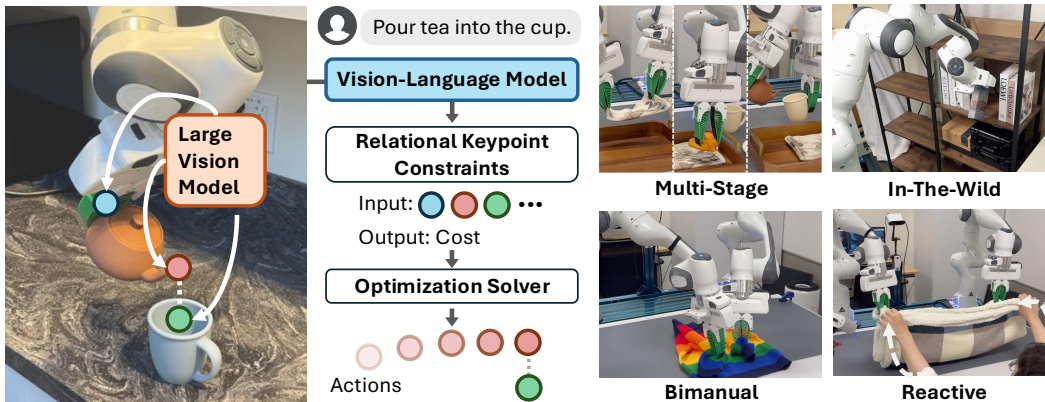


Figure 1: Relational Keypoint Constraints (ReKep) specify diverse manipulation behaviors as an optimizable spatio-temporal series of constraint functions operating on semantic keypoints. In the pouring task, one ReKep first constrains the grasping location at the handle of the teapot (*blue*). A subsequent ReKep pulls the teapot spout (*red*) towards the top of the cup opening (*green*) while another ReKep constrains the desired rotation of the teapot by associating the vector formed by the handle (*blue*) and the spout (*red*).

21 1 Introduction

22 Robotic manipulation involves intricate interactions with objects in the environment, which can of-
23 ten be expressed as constraints in both spatial and temporal domains. Consider the task of pouring
24 tea into a cup in Fig. 1: the robot must grasp *at the handle*, keep the cup *upright* while transporting
25 it, *align* the spout with the target container, and then tilt the cup *at the correct angle* to pour. Here,
26 the constraints encode not only the intermediate sub-goals (e.g., align the spout) but also the transi-
27 tioning behaviors (e.g., keep the cup *upright* in transportation), which collectively dictate the spatial,
28 timing, and other combinatorial requirements of the robot’s actions in relation to the environment.

29 However, effectively formulating these constraints for a large variety of real-world tasks presents
30 significant challenges. While representing constraints using relative poses between robots and ob-
31 jects is a direct and widely-used approach [1], rigid-body transformations do not depict geometric
32 details, require obtaining object models *a priori*, and cannot work on deformable objects. On the
33 other hand, data-driven approaches enable learning constraints directly in visual space [2, 3]. While
34 more flexible, it remains unclear how to effectively collect training data as the number of constraints
35 grows combinatorially in terms of objects and tasks. Therefore, we ask the question: how can we
36 represent constraints in manipulation that are 1) widely applicable: adaptable to tasks that require
37 multi-stage, in-the-wild, bimanual, and reactive behaviors, 2) scalably obtainable: have the potential
38 to be fully automated through the advances in foundation models, and 3) real-time optimizable: can
39 be efficiently solved by off-the-shelf solvers to produce complex manipulation behaviors?

40 In this work, we propose **Relational Keypoint Constraints (ReKep)**. Specifically, ReKep repre-
41 sents constraints as Python functions that map a set of keypoints to a numerical cost, where each
42 keypoint is a task-specific and semantically meaningful 3D point in the scene. Each function is
43 composed of (potentially nonlinear) arithmetic operations on the keypoints and encodes a desired
44 “relation” between them, where the keypoints may belong to different entities in the environment,
45 such as the robot arms, object parts, and other agents. While each keypoint only consists of its 3D
46 Cartesian coordinates in the world frame, multiple keypoints can collectively specify lines, surfaces,
47 and/or 3D rotations if rigidity between keypoints is enforced. We study ReKep in the context of the
48 sequential manipulation problem, where each task involves multiple stages that have spatio-temporal
49 dependencies (e.g., “grasping”, “aligning”, and “pouring” in the aforementioned example).

50 While constraints are typically defined manually per task [4], we demonstrate the specific form
51 of ReKep possesses a unique advantage in that they can be automated by pre-trained large vi-
52 sion models (LVM) [5] and vision-language models (VLM) [6], enabling *in-the-wild* specification
53 of ReKep from RGB-D observations and free-form language instructions. Specifically, we leverage
54 LVM to propose fine-grained and semantically meaningful keypoints in the scene and VLM to write
55 the constraints as Python functions from visual input overlaid with proposed keypoints. This process
56 can be interpreted as grounding fine-grained spatial relations, often those not easily specified with
57 natural language, in an output modality supported by VLM (code) using visual referral expressions.

58 With the generated constraints, off-the-shelf solvers can be used to produce robot actions by re-
59 evaluating the constraints based on tracked keypoints. Inspired by [7], we employ a hierarchical
60 optimization procedure to first solve a set of waypoints as sub-goals (represented as $SE(3)$ end-
61 effector poses) and then solve the receding-horizon control problem to obtain a dense sequence of
62 actions to achieve each sub-goal. With appropriate instantiation of the problem, we demonstrate that
63 it can be reliably solved at approximately 10 Hz for the tasks considered in this work.

64 Our contributions are summarized as follows: 1) We formulate manipulation tasks as a hierarchical
65 optimization problem with Relational Keypoint Constraints; 2) We devise a pipeline to automatically
66 specify keypoints and constraints using large vision models and vision-language models; 3) We
67 present system implementations on two real-robot platforms that take as input a language instruction
68 and RGB-D observations, and produce multi-stage, in-the-wild, bimanual, and reactive behaviors for
69 a large variety of manipulation tasks, all without task-specific data or environment models.

70 2 Related Works

71 **Structural Representations for Manipulation.** Structural representations determine the orchestra-
72 tion of different modules in a manipulation system and yield different implications on the capabili-
73 ties, assumptions, efficiency, and effectiveness of the system. Rigid-body poses are most commonly
74 used given well-understood rigid-body motions in free space and their efficiency at modeling long-
75 range dependencies of objects [1, 8–18]. However, since it often requires both geometry and dynam-
76 ics of environment to be modeled beforehand, various works have studied structural representations
77 using data-driven methods, such as learning object-centric representation [19–34], particle-based
78 dynamics [35–41], and keypoints or descriptors [3, 4, 42–54]. Among them, keypoints have shown

79 great promises given their interpretability, efficiency, generalization to instance variations [4], and
80 ability to model both rigid bodies and deformable objects. However, manual annotation is required
81 per task, thus lacking scalability in open-world settings, which we aim to address in this work.

82 **Constrained Optimization in Manipulation.** Constraints are often used to impose desired be-
83 haviors on robots. Motion planning algorithms use geometric constraints to compute feasible tra-
84 jectories that avoid obstacles and achieve goals [55–60]. Contact constraints can be used to plan
85 forceful or contact-rich behaviors [61–74]. For sequential manipulation tasks, task and motion plan-
86 ning (TAMP) [1, 8, 13] is a widely used framework, often formulated as constraint satisfaction
87 problems [11, 75–80] with continuous geometric problems as subroutines. Logic-Geometric Pro-
88 gramming [81–86] alternatively formulates a nonlinear constrained program over the entire state
89 trajectory, taking into account both logical and geometric constraints. Constraints can either be
90 manually written or learned from data in the form of manifolds [87], feasibility model [86, 88], or
91 signed-distance fields [2, 89]. Inspired by [7], we formulate sequential manipulation tasks as an
92 integrated *continuous* mathematical program that is repeatedly solved in a receding-horizon fashion,
93 with the key difference being that the constraints are synthesized by foundation models.

94 **Foundation Models for Robotics.** Leveraging foundation models for robotics is an active area
95 of research. We refer readers to [90–93] for overview and recent applications. Here, we focus on
96 VLMs that are capable of incorporating visual inputs [6, 94–98] for manipulation. However, despite
97 showing promises for open-world planning and goal specification [99–113], the caption-guided pre-
98 training scheme of VLMs often limits visual details that can be retained about the image [114–117].
99 Self-supervised vision models (e.g., DINO [5, 118]), on the other hand, provide fine-grained pixel-
100 level features useful for various vision and robotic tasks [31, 119–124], but there lack effective ways
101 for interpreting open-world semantics pivotal for cross-task generalization. In this work, we leverage
102 their complementary strengths by using DINOv2 [5] for fine-grained keypoint proposal and using
103 GPT-4o [6] for its visual reasoning capability in a supported output modality (code). Similar forms
104 of visual prompting techniques are also explored in concurrent works [99–101, 112, 125]. In this
105 work, we demonstrate ReKep possesses the unique advantages of performing challenging 6-12 DoF
106 tasks, integrated high-level reasoning for reactive replanning, high-frequency closed-loop execution,
107 and generating black-box constraints via visual prompting. More discussion in Appendix A.10.

108 3 Method

109 Herein we discuss: **(1)** What are Relational Keypoint Constraints (Sec. 3.1)? **(2)** How to formulate
110 manipulation as a constrained optimization problem with ReKep (Sec. 3.2)? **(3)** What is our algo-
111 rithmic instantiation that can efficiently solve the optimization in real-time (Sec. 3.3)? **(4)** How to
112 automatically obtain ReKep from RGB-D observations and language instructions (Sec. 3.4)?

113 3.1 Relational Keypoint Constraints (ReKep)

114 Herein we define a single instance of ReKep. For clarity, we assume that a set of K keypoints have
115 been specified (discussed later in Sec. 3.4). Concretely, each keypoint $k_i \in \mathbb{R}^3$ refers to a 3D point
116 on the scene surface with Cartesian coordinates, which is dependent on the task semantics and the
117 environment (e.g., grasp point on the handle, spout).

118 A single instance of ReKep is a function $f : \mathbb{R}^{K \times 3} \rightarrow \mathbb{R}$ that maps an array of keypoints, denoted
119 as \mathbf{k} , to an unbounded cost, where $f(\mathbf{k}) \leq 0$ indicates the constraint is satisfied. The function f is
120 implemented as a stateless Python function, containing NumPy [126] operations on keypoints, which
121 may be nonlinear and nonconvex. In essence, one instance of ReKep encodes one desired spatial
122 relation between keypoints, which may belong to robot arm(s), object parts, and other agents.

123 However, a manipulation task typically involves multiple spatial relations and may have multiple
124 temporally dependent stages where each stage entails different spatial relations. To this end, we
125 decompose a task into N stages and use ReKep to specify two kinds of constraints for each stage
126 $i \in \{1, \dots, N\}$: a set of sub-goal constraints $\mathcal{C}_{\text{sub-goal}}^{(i)} = \{f_{\text{sub-goal},1}^{(i)}(\mathbf{k}), \dots, f_{\text{sub-goal},n}^{(i)}(\mathbf{k})\}$ and a

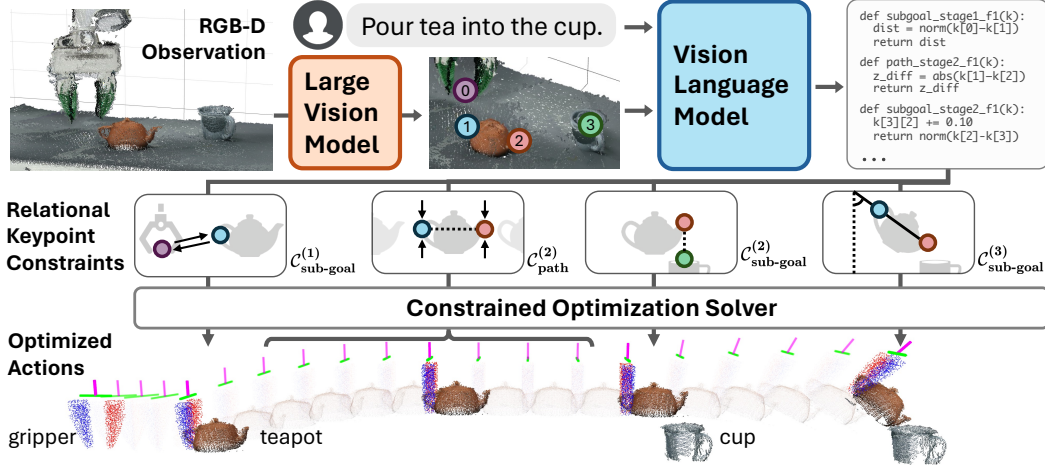


Figure 2: Overview of ReKep. Given RGB-D observation and free-form language instruction, DINOv2 [5] is used to propose keypoint candidates on fine-grained meaningful regions in the scene. The image overlaid with keypoints and the instruction are fed into GPT-4o [6] to generate a series of ReKep constraints as python programs that specify desired relations between keypoints at different stages of the task ($\mathcal{C}_{\text{sub-goal}}^{(i)}$) and any requirement on the transitioning behaviors ($\mathcal{C}_{\text{path}}^i$). Finally, a constrained optimization solver is used to obtain a dense sequence of end-effector actions in $SE(3)$, subject to the generated constraints.

127 set of path constraints $\mathcal{C}_{\text{path}}^{(i)} = \{f_{\text{path},1}^{(i)}(\mathbf{k}), \dots, f_{\text{path},m}^{(i)}(\mathbf{k})\}$, where $f_{\text{sub-goal}}^{(i)}$ encodes one keypoint
 128 relation to be achieved *at the end of stage i* , and $f_{\text{path}}^{(i)}$ encodes one keypoint relation to be satisfied
 129 for every state *within stage i* . Consider the pouring task in Fig. 2, which consists of three stages:
 130 grasp, align, and pour. The stage-1 sub-goal constraint pulls the end-effector towards the teapot
 131 handle. Then stage-2 sub-goal constraint specifies that the teapot spout needs to be on top of the
 132 cup opening. Additionally, stage-2 path constraint ensures the teapot stays upright to avoid spillage
 133 when transported. Finally, the stage-3 sub-goal constraint specifies the desired pouring angle.

134 3.2 Manipulation Tasks as Constrained Optimization with ReKep

135 Using ReKep as a general tool to represent constraints, we adopt the formulation in [7] and show
 136 how a manipulation task can be formulated as a constrained optimization problem involving $\mathcal{C}_{\text{sub-goal}}^{(i)}$
 137 and $\mathcal{C}_{\text{path}}^{(i)}$. We denote the end-effector pose as $\mathbf{e} \in SE(3)$. To perform the manipulation task, we aim
 138 to obtain the overall discrete-time trajectory $\mathbf{e}_{1:T}$ by formulating the control problem as follows:

$$\arg \min_{\mathbf{e}_{1:T}, g_{1:N}} \sum_{i=1}^N \left[\lambda_{\text{sub-goal}}^{(i)}(\mathbf{e}_{g_i}) + \sum_{t=g_{i-1}}^{g_i} \lambda_{\text{path}}^{(i)}(\mathbf{e}_t) \right] \text{ s.t. } \begin{cases} \mathbf{e}_1 = \mathbf{e}_{\text{init}}, g_0 = 1, 0 < g_i < g_{i+1} \\ f(\mathbf{k}_{g_i}) \leq 0, \forall f \in \mathcal{C}_{\text{sub-goal}}^{(i)} \\ f(\mathbf{k}_t) \leq 0, \forall f \in \mathcal{C}_{\text{path}}^{(i)}, t = g_{i-1}, \dots, g_i \\ \mathbf{k}_{t+1} = h(\mathbf{k}_t, \mathbf{e}_t), t = 1, \dots, T-1 \end{cases} \quad (1)$$

139 where \mathbf{e}_t denotes end-effector pose at time t , $g_i \in \{1, \dots, T\}$ are the timings of the transition from
 140 stage i to $i+1$ which are also auxiliary decision variables, \mathbf{k}_t is the array of keypoint positions at
 141 time t , h is a forward model of keypoints, and $\lambda_{\text{sub-goal}}^{(i)}$ and $\lambda_{\text{path}}^{(i)}$ are auxiliary cost functions (e.g.,
 142 collision avoidance) for the sub-goal and path problems respectively. Namely, for each stage i , the
 143 optimization shall find an end-effector pose as next sub-goal, along with its timing, and a sequence of
 144 poses $\mathbf{e}_{g_{i-1}:g_i}$ that achieves the sub-goal, subject to the given set of ReKep constraints and auxiliary
 145 costs. This formulation can be considered as direct shooting in trajectory optimization [127].
 146

147 3.3 Decomposition and Algorithmic Instantiation

148 To solve Eq. 1 in real-time, we employ a decomposition of the full problem and only optimize
 149 for the immediate next sub-goal and the corresponding path to reach the sub-goal (pseudo-code

150 in Algorithm 1). All optimization problems are implemented and solved using SciPy [128] with
 151 decision variables normalized to $[0, 1]$. They are initially solved with Dual Annealing [129] with
 152 SLSQP [130] as local optimizer (around 1 second) and subsequently solved with only local opti-
 153 mizer based on the previous solution at approximately 10 Hz.

154 **The Sub-Goal Problem:** We first solve the sub-goal problem to obtain \mathbf{e}_{g_i} for the current stage i :

$$\arg \min_{\mathbf{e}_{g_i}} \lambda_{\text{sub-goal}}^{(i)}(\mathbf{e}_{g_i}) \quad \text{s.t.} \quad f(\mathbf{k}_{g_i}) \leq 0, \quad \forall f \in \mathcal{C}_{\text{sub-goal}}^{(i)} \quad (2)$$

155 where $\lambda_{\text{sub-goal}}$ subsumes auxiliary control costs: scene collision avoidance, reachability, pose regu-
 156 larization, solution consistency, and self-collision for bimanual setup (details in A.8). Namely, Eq. 2
 157 attempts to find a sub-goal that satisfies $\mathcal{C}_{\text{sub-goal}}^i$ while minimizing the auxiliary costs. If a stage is
 158 concerned with grasping, a grasp metric is also included. In this work, we use AnyGrasp [131]¹.

159 **The Path Problem:** After obtaining sub-goal \mathbf{e}_{g_i} , we solve for a trajectory $\mathbf{e}_{t:g_i}$ starting from current
 160 end-effector pose \mathbf{e}_t to the sub-goal \mathbf{e}_{g_i} :

$$\arg \min_{\mathbf{e}_{t:g_i, g_i}} \lambda_{\text{path}}^{(i)}(\mathbf{e}_{t:g_i}) \quad \text{s.t.} \quad f(\mathbf{k}_{\hat{t}}) \leq 0, \quad \forall f \in \mathcal{C}_{\text{path}}^{(i)}, \quad \hat{t} = t, \dots, g_i \quad (3)$$

161 where λ_{path} subsumes the following auxiliary control costs: scene collision avoidance, reachability,
 162 path length, solution consistency, and self-collision in the case of bimanual setup (details in A.9). If
 163 the distance to the sub-goal \mathbf{e}_{g_i} is within a small tolerance ϵ , we progress to the next stage $i + 1$.

164 **Backtracking:** Although the sub-problems can be solved at a real-time frequency to react to external
 165 disturbances within a stage, it is imperative that the system can replan across stages if any sub-goal
 166 constraint from the last stage no longer holds (e.g., cup taken out of the gripper in the pouring task).
 167 Specifically, in every control loop, we check for violation of $\mathcal{C}_{\text{path}}^{(i)}$. If one is found, we iteratively
 168 backtrack to a previous stage j such that $\mathcal{C}_{\text{path}}^{(j)}$ is satisfied.

169 **Forward Models for Keypoints:** To solve Eq. 2 and Eq. 3, one must utilize a forward model h
 170 that estimates $\Delta \mathbf{k}$ from $\Delta \mathbf{e}$ in the optimization process. As in prior work [4], we make the rigidity
 171 assumption between the end-effector and the ‘‘grasped keypoints’’ (a rigid group of keypoints that
 172 belong to the same object or part; obtained from the segmentation model as described in Sec. 3.4.).
 173 Namely, given a change in the end-effector pose $\Delta \mathbf{e}$, we can calculate the change in keypoint posi-
 174 tions by applying the same relative rigid transformation $\mathbf{k}'[\text{grasped}] = \mathbf{T}_{\Delta \mathbf{e}} \cdot \mathbf{k}[\text{grasped}]$, while
 175 assuming other keypoints stay static. We note that this is a ‘‘local’’ assumption in that it is only as-
 176 sumed to hold for the short horizon (0.1s) that the problem is solved. Actual keypoint positions are
 177 tracked using visual input at 20 Hz and used in every new problem. For more challenging scenarios
 178 (e.g., dynamic or contact-rich tasks), a learned or physics-based model may be used.

179 3.4 Keypoint Proposal and ReKep Generation

180 To enable the system to perform tasks in-the-wild given a free-form task instruction, we devise a
 181 pipeline using large vision models and vision-language models for keypoint proposal and ReKep
 182 generation, which are respectively discussed as follows:

183 **Keypoint Proposal:** Given an RGB image $\mathbb{R}^{h \times w \times 3}$, we first extract the patch-wise features $\mathbf{F}_{\text{patch}} \in$
 184 $\mathbb{R}^{h' \times w' \times d}$ from DINOv2 [5]. Then we perform bilinear interpolation to upsample the features to
 185 the original image size, $\mathbf{F}_{\text{interp}} \in \mathbb{R}^{h \times w \times d}$. To ensure the proposal covers all relevant objects in
 186 the scene, we extract all masks $\mathbf{M} = \{\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_n\}$ in the scene using Segment Anything
 187 (SAM) [132]. For each mask j , we cluster the masked features $\mathbf{F}_{\text{interp}}[\mathbf{m}_j]$ using k -means with
 188 $k = 5$ with a cosine-similarity metric. The centroids of the clusters are used as keypoint candidates,
 189 which are projected to a world coordinate \mathbb{R}^3 using a calibrated RGB-D camera. Candidates that are
 190 within 8cm of others are filtered out. Overall, we find that this procedure is adept at identifying a
 191 large percentage of fine-grained and semantically meaningful regions of objects.

¹Since AnyGrasp is a grasp detector instead of a metric and is computationally expensive to run in optimiza-
 tion loops, we always return the grasp closest to a specified ‘‘grasp keypoint’’ by exploiting the fact that ReKep
 related to grasping always associates a dummy keypoint on the end-effector and one actual keypoint.

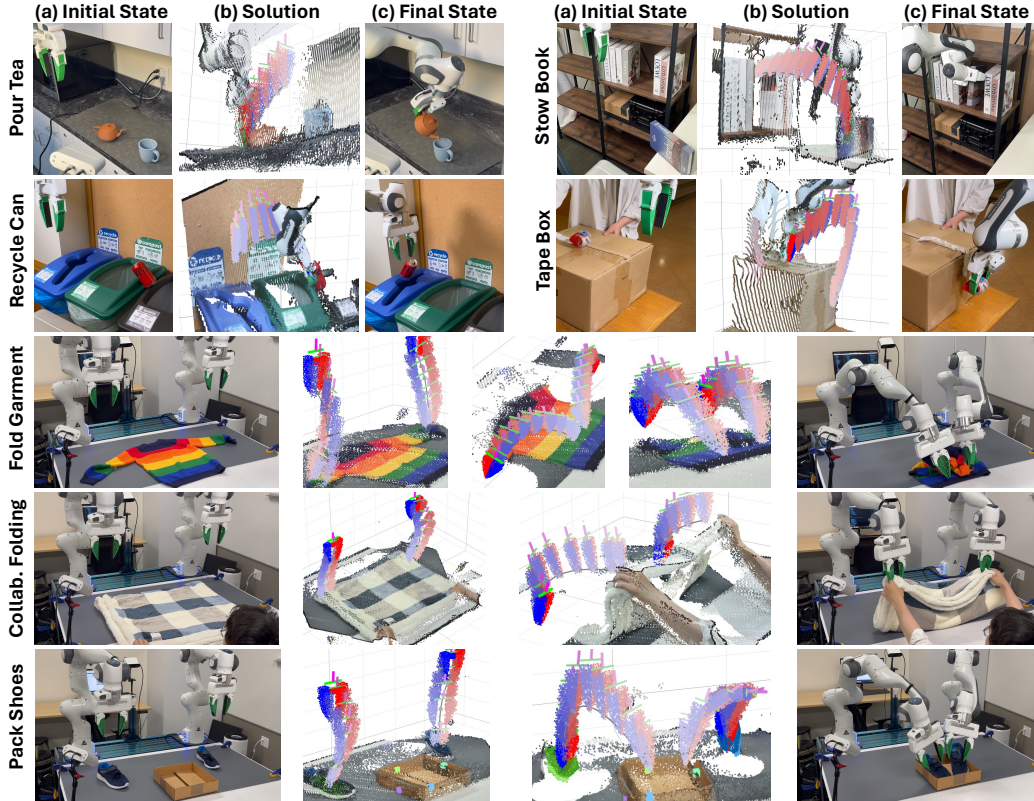


Figure 3: Experiment tasks and visualization of optimization results. Seven tasks are designed to validate different aspects of our system, including in-the-wild specification with commonsense knowledge, multi-stage tasks with spatio-temporal dependencies, bimanual coordination with geometric awareness, and reactivity when collaborating with humans and under disturbances.

192 **ReKep Generation:** After obtaining the keypoint candidates, we overlay them on the original RGB
 193 image with numerical marks. Coupled with the language instruction of the task, we then use visual
 194 prompting to query GPT-4o [6] to generate the number of required stages and the corresponding sub-
 195 goal constraints $\mathcal{C}_{\text{sub-goal}}^{(i)}$ and path constraints $\mathcal{C}_{\text{path}}^{(i)}$ for each stage i (prompts are in A.6). Notably,
 196 the functions do not directly manipulate the numerical values of the keypoint positions. Rather, we
 197 exploit the strength of VLM to specify spatial relations as *arithmetic operations*, such as L2 distance
 198 or dot product between keypoints, that are only instantiated when invoked with actual keypoint posi-
 199 tions tracked by a specialized 3D tracker. Furthermore, an important advantage of using arithmetic
 200 operations on a set of keypoint positions is that it can specify 3D rotations in full $SO(3)$ when suffi-
 201 cient points are provided and rigidity between relevant points is enforced, but this is done only when
 202 needed depending on task semantic. This enables VLM to reason about 3D rotations with arithmetic
 203 operations in *3D Cartesian space*, effectively circumventing the need for dealing with alternative
 204 3D rotation representation and the need for performing numerical computation.

205 4 Experiments

206 We aim to answer the following research questions: (1) How well does our framework automatically
 207 formulate and synthesize manipulation behaviors (Sec. 4.1)? (2) Can our system generalize to novel
 208 objects and manipulation strategies (Sec. 4.2)? (3) How do the individual components contribute to
 209 the failure cases of the system (Sec. 4.3)? We validate ReKep on two real robot platforms: a wheeled
 210 single-arm platform, and a stationary dual-arm platform (Figure. 3). Additional implementation
 211 details can be found in Appendix, including keypoint proposal (A.5), VLM querying (A.6), point
 212 trackers (A.7), sub-goal solver (A.8), and path solver (A.9).

213 **4.1 In-the-Wild and Bimanual Manipulation with ReKep**

214 **Tasks.** We purposefully select a set of tasks (shown
 215 in Fig. 3) with the goal of examining the multi-stage
 216 (**m**), in-the-wild (**w**), bimanual (**b**), and reactive (**r**) be-
 217 haviors of the system. The tasks and their features are
 218 *Pour Tea* (**m, w, r**), *Stow Book* (**w**), *Recycle Can* (**w**),
 219 *Tape Box* (**w, r**), *Fold Garment* (**b**), *Pack Shoes* (**b**),
 220 and *Collaborative Folding* (**b, r**). We further evaluate
 221 three of the tasks under external disturbances (denoted
 222 as “Dist.”) by changing poses of task objects during
 223 execution.

224 **Metric and Baselines.** Each setting has 10 trials, in
 225 which object poses are randomized. Task success rate is
 226 reported in Tab. 1. We compare against VoxPoser [106]
 227 as a baseline. We evaluate two variants of the system,
 228 “Auto” which uses foundation models for automating
 229 generation of ReKep, and “Annotated (Annot.)” which
 230 uses human-annotated ReKep.

231 **Results.** Overall the system demonstrates promising
 232 capabilities at formulating correct constraints and ex-
 233 ecuting them in unstructured environments, despite no
 234 task-specific data or environment model are provided.
 235 Notably, ReKep can effectively handle core challenges
 236 of each task. For example, it can formulate correct tem-
 237 poral dependency in multi-stage tasks (e.g., spout needs
 238 to be aligned with the cup before pouring), leverage commonsense knowledge (e.g., coke cans
 239 should be recycled), and construct coordination behaviors in both bimanual settings (e.g., fold-
 240 ing left sleeve and right sleeve simultaneously) and human-robot collaboration setting (e.g., folding
 241 a large blanket by aligning the four corners together with human). Coupled with an optimization
 242 framework, it can also generate kinematically challenging behaviors in confined spaces in the *Stow*
 243 *Book* task and find a feasible solution that densely fits two shoes within a small volume in the *Pack*
 244 *Shoes* task. Since the keypoints are tracked at a high frequency, the system can also react to external
 245 disturbances and replan both within stage and across stages. However, despite promising results, we
 246 also observe that the generated constraints are not always fully correct but expect to see improvement
 247 given the rapid advancement in these pre-trained models.

248 **4.2 Generalization in Manipulation Strategies**

249 **Tasks.** We systematically evaluate how novel manipulation strategies can be formulated by focusing
 250 on a single task, garment folding, but with 8 unique categories of garments, each demanding a unique
 251 way of folding and requiring both geometrical and commonsense reasoning. Evaluation is done on
 252 the bimanual platform, presenting additional challenges in bimanual coordination.

253 **Metric.** We use GPT-4o with a prompt containing only generic instructions with no in-context
 254 examples. “Strategy Success” measures whether generated ReKep is feasible, which tests both the
 255 keypoint proposal module and the VLM, and “Execution Success” measures system success rate
 256 given feasible strategies for each clothing. Each is measured with 10 trials.

257 **Results.** Interestingly, we observe drastically different strategies across categories, many of which
 258 are aligned with how humans might fold each garment. For example, it can recognize that two
 259 sleeves often are folded together, prior to fully folding the clothes. In cases where using two arms
 260 is unnecessary, akin to how humans fold clothes, only one arm is being used. However, we do
 261 observe that the VLM may miss certain steps to complete the folding as the operator expected, but
 262 we recognize that this is inherently an open-ended problem often based on one’s preferences.

Task	VoxPoser	ReKep	
		Auto	Annot.
Mobile Arm			
Pour Tea	0/10	3/10	8/10
Recycle Can	3/10	6/10	8/10
Stow Book	0/10	3/10	6/10
Tape Box	4/10	7/10	8/10
Dual-Arm			
Fold Garment	0/10	5/10	6/10
Pack Shoes	0/10	3/10	5/10
Collab. Folding	0/10	4/10	7/10
Total (%)	10.0%	44.3%	68.6%

Table 1: Success rate on the wheeled single-arm and the stationary bimanual platforms.

Task (Dist.)	VoxPoser	ReKep	
		Auto	Annot.
Pour Tea	0/10	2/10	4/10
Tape Box	2/10	3/10	5/10
Collab. Folding	0/10	3/10	5/10
Total (%)	6.7%	26.7%	46.7%

Table 2: Success rate under external disturbances across both robot platforms.









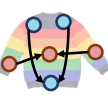
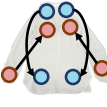


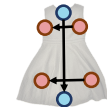
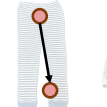
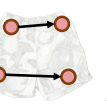
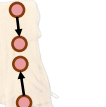
	sweater	shirt	hoodie	vest	dress	pants	shorts	scarf	Total
Garment									
ReKep									
Strategy Success	6/10	4/10	4/10	6/10	3/10	7/10	7/10	5/10	52.5%
Execution Success	6/10	5/10	6/10	9/10	7/10	8/10	9/10	9/10	73.8%

Figure 4: Novel *bimanual* strategies of ReKep for folding different categories of garments and their success rates. Since ReKep in this task always associates two points at a time, two keypoints are connected by an arrow if they need to be aligned. The coloring of the keypoints denotes the order. In the sweater task, two sleeves are first folded simultaneously with two arms, and then the two arms grasp the crew neck to align to the bottom.

263 4.3 System Error Breakdown

264 The modular design of the framework entails an advantage for analyzing system errors due to its interpretability. In this section, we
 265 perform an empirical investigation by manually inspecting the failure cases of the experiments reported in Tab. 1, which is then used to
 266 calculate the likelihood of a module causing an error while accounting for their temporal dependencies in the pipeline. Results are reported
 267 in Fig. 5. Among the different modules, the point tracker produces the largest portion of errors, as frequent and intermittent
 268 occlusion poses significant challenges for accurate tracking. Key-
 269 point proposal and VLM also produce considerable portions of errors, where common cases include
 270 the proposal module missing certain keypoints and the VLM referring to incorrect keypoints. The
 271 optimization module, on the other hand, does not contribute as much to the failures despite given
 272 limited time budget, since there often exist many possible solutions for each problem. Other mod-
 273 ules, such as segmentation, 3D reconstruction, and low-level controller, also contribute to some
 274 failure cases, but they are relatively insignificant compared to other modules.

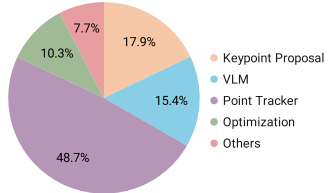


Figure 5: Error breakdown of the system modules.

279 5 Conclusion & Limitations

280 In this work, we presented Relational Keypoint Constraints (ReKep), a structural task representa-
 281 tion using constraints that operates on semantic keypoints to specify desired relations between robot
 282 arms, object (parts), and other agents in the environment. Coupled with point trackers, we demon-
 283 strate that ReKep constraints can be repeatedly and efficiently solved in a hierarchical optimization
 284 framework to act as a closed-loop policy that runs at a real-time frequency. We also demonstrate
 285 the unique advantage of ReKep in that it can be automatically synthesized by large vision models
 286 and vision-language models. Results are shown on two robot platforms and on a variety of tasks
 287 featuring multi-stage, in-the-wild, bimanual, and reactive behaviors, all without task-specific data,
 288 additional training, or environment models.

289 Despite the promises, several limitations remained. First, the optimization framework relies on a
 290 forward model of keypoints based on rigidity assumption, albeit a high-frequency feedback loop that
 291 relaxes the accuracy requirement of the model. Second, ReKep relies on accurate point tracking to
 292 correctly optimize actions in closed-loop, which is itself a challenging 3D vision task due to heavy
 293 intermittent occlusions. Lastly, the current formulation assumes a fixed sequence of stages (i.e.,
 294 skeletons) for each task. Replanning with different skeletons requires running keypoint proposal
 295 and VLM at a high-frequency, which poses considerable computational challenges. An extended
 296 discussion can be found in Appendix A.11.

References

- 297
- 298 [1] L. P. Kaelbling and T. Lozano-Pérez. Hierarchical planning in the now. In *Workshops at the*
299 *Twenty-Fourth AAAI Conference on Artificial Intelligence*, 2010.
- 300 [2] D. Driess, J.-S. Ha, M. Toussaint, and R. Tedrake. Learning models as functionals of signed-
301 distance fields for manipulation planning. In *Conference on robot learning*, pages 245–255.
302 PMLR, 2022.
- 303 [3] A. Simeonov, Y. Du, A. Tagliasacchi, J. B. Tenenbaum, A. Rodriguez, P. Agrawal, and V. Sitz-
304 mann. Neural descriptor fields: Se (3)-equivariant object representations for manipulation.
305 In *2022 International Conference on Robotics and Automation (ICRA)*, pages 6394–6400.
306 IEEE, 2022.
- 307 [4] L. Manuelli, W. Gao, P. Florence, and R. Tedrake. kpm: Keypoint affordances for category-
308 level robotic manipulation. In *The International Symposium of Robotics Research*, pages
309 132–157. Springer, 2019.
- 310 [5] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez,
311 D. Haziza, F. Massa, A. El-Nouby, et al. Dinov2: Learning robust visual features without
312 supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 313 [6] OpenAI. Gpt-4 technical report. *arXiv*, 2023.
- 314 [7] M. Toussaint, J. Harris, J.-S. Ha, D. Driess, and W. Hönig. Sequence-of-constraints mpc:
315 Reactive timing-optimal control of sequential manipulation. In *2022 IEEE/RSJ International*
316 *Conference on Intelligent Robots and Systems (IROS)*, pages 13753–13760. IEEE, 2022.
- 317 [8] L. P. Kaelbling and T. Lozano-Pérez. Integrated task and motion planning in belief space.
318 *The International Journal of Robotics Research*, 32(9-10):1194–1227, 2013.
- 319 [9] S. Srivastava, E. Fang, L. Riano, R. Chitnis, S. Russell, and P. Abbeel. Combined task and
320 motion planning through an extensible planner-independent interface layer. In *2014 IEEE*
321 *international conference on robotics and automation (ICRA)*, 2014.
- 322 [10] A. Byravan and D. Fox. Se3-nets: Learning rigid body motion using deep neural networks.
323 In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 173–180.
324 IEEE, 2017.
- 325 [11] N. T. Dantam, Z. K. Kingston, S. Chaudhuri, and L. E. Kavraki. An incremental constraint-
326 based framework for task and motion planning. *The International Journal of Robotics Re-*
327 *search*, 37(10):1134–1151, 2018.
- 328 [12] T. Migimatsu and J. Bohg. Object-centric task and motion planning in dynamic environments.
329 *IEEE Robotics and Automation Letters*, 5(2):844–851, 2020.
- 330 [13] C. R. Garrett, R. Chitnis, R. Holladay, B. Kim, T. Silver, L. P. Kaelbling, and T. Lozano-Pérez.
331 Integrated task and motion planning. *Annual review of control, robotics, and autonomous*
332 *systems*, 4:265–293, 2021.
- 333 [14] A. Curtis, X. Fang, L. P. Kaelbling, T. Lozano-Pérez, and C. R. Garrett. Long-horizon ma-
334 nipulation of unknown objects via task and motion planning with estimated affordances. In
335 *2022 International Conference on Robotics and Automation (ICRA)*, pages 1940–1946. IEEE,
336 2022.
- 337 [15] Y. Labbé, L. Manuelli, A. Mousavian, S. Tyree, S. Birchfield, J. Tremblay, J. Carpentier,
338 M. Aubry, D. Fox, and J. Sivic. Megapose: 6d pose estimation of novel objects via render &
339 compare. In *Proceedings of the 6th Conference on Robot Learning (CoRL)*, 2022.

- 340 [16] S. Tyree, J. Tremblay, T. To, J. Cheng, T. Mosier, J. Smith, and S. Birchfield. 6-dof pose esti-
341 mation of household objects for robotic manipulation: An accessible dataset and benchmark.
342 In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages
343 13081–13088. IEEE, 2022.
- 344 [17] C. Pan, B. Okorn, H. Zhang, B. Eisner, and D. Held. Tax-pose: Task-specific cross-pose
345 estimation for robot manipulation. In *Conference on Robot Learning*, pages 1783–1792.
346 PMLR, 2023.
- 347 [18] B. Wen, W. Yang, J. Kautz, and S. Birchfield. Foundationpose: Unified 6d pose estimation
348 and tracking of novel objects. *arXiv preprint arXiv:2312.08344*, 2023.
- 349 [19] I. Lenz, R. A. Knepper, and A. Saxena. Deepmpc: Learning deep latent features for model
350 predictive control. In *Robotics: Science and Systems*, volume 10. Rome, Italy, 2015.
- 351 [20] M. B. Chang, T. Ullman, A. Torralba, and J. B. Tenenbaum. A compositional object-based
352 approach to learning physical dynamics. *arXiv preprint arXiv:1612.00341*, 2016.
- 353 [21] P. Battaglia, R. Pascanu, M. Lai, D. Jimenez Rezende, et al. Interaction networks for learning
354 about objects, relations and physics. *Advances in neural information processing systems*, 29,
355 2016.
- 356 [22] A. Sanchez-Gonzalez, N. Heess, J. T. Springenberg, J. Merel, M. Riedmiller, R. Hadsell,
357 and P. Battaglia. Graph networks as learnable physics engines for inference and control. In
358 *International Conference on Machine Learning*, pages 4470–4479. PMLR, 2018.
- 359 [23] E. Jang, C. Devin, V. Vanhoucke, and S. Levine. Grasp2vec: Learning object representations
360 from self-supervised grasping. *arXiv preprint arXiv:1811.06964*, 2018.
- 361 [24] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield. Deep ob-
362 ject pose estimation for semantic robotic grasping of household objects. *arXiv preprint*
363 *arXiv:1809.10790*, 2018.
- 364 [25] Z. Xu, J. Wu, A. Zeng, J. B. Tenenbaum, and S. Song. Densephysnet: Learning dense physical
365 object representations via multi-step dynamic interactions. *arXiv preprint arXiv:1906.03853*,
366 2019.
- 367 [26] J. Mao, C. Gan, P. Kohli, J. B. Tenenbaum, and J. Wu. The neuro-symbolic concept
368 learner: Interpreting scenes, words, and sentences from natural supervision. *arXiv preprint*
369 *arXiv:1904.12584*, 2019.
- 370 [27] C. P. Burgess, L. Matthey, N. Watters, R. Kabra, I. Higgins, M. Botvinick, and A. Ler-
371 chner. Monet: Unsupervised scene decomposition and representation. *arXiv preprint*
372 *arXiv:1901.11390*, 2019.
- 373 [28] L. Hewing, K. P. Wabersich, M. Menner, and M. N. Zeilinger. Learning-based model pre-
374 dictive control: Toward safe learning in control. *Annual Review of Control, Robotics, and*
375 *Autonomous Systems*, 3:269–296, 2020.
- 376 [29] F. Locatello, D. Weissenborn, T. Unterthiner, A. Mahendran, G. Heigold, J. Uszkoreit,
377 A. Dosovitskiy, and T. Kipf. Object-centric learning with slot attention. *Advances in neural*
378 *information processing systems*, 33:11525–11538, 2020.
- 379 [30] N. Heravi, A. Wahid, C. Lynch, P. Florence, T. Armstrong, J. Tompson, P. Sermanet, J. Bohg,
380 and D. Dwibedi. Visuomotor control in multi-object scenes using object-aware representa-
381 tions. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages
382 9515–9522. IEEE, 2023.
- 383 [31] Y. Zhu, Z. Jiang, P. Stone, and Y. Zhu. Learning generalizable manipulation policies with
384 object-centric 3d representations. *arXiv preprint arXiv:2310.14386*, 2023.

- 385 [32] W. Yuan, C. Paxton, K. Desingh, and D. Fox. Soronet: Spatial object-centric representations
386 for sequential manipulation. In *Conference on Robot Learning*, pages 148–157. PMLR, 2022.
- 387 [33] S. Cheng, C. Garrett, A. Mandlekar, and D. Xu. Nod-tamp: Multi-step manipulation planning
388 with neural object descriptors. *arXiv preprint arXiv:2311.01530*, 2023.
- 389 [34] J. Hsu, J. Mao, J. Tenenbaum, and J. Wu. What’s left? concept grounding with logic-enhanced
390 foundation models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 391 [35] Y. Li, J. Wu, R. Tedrake, J. B. Tenenbaum, and A. Torralba. Learning particle dynamics for
392 manipulating rigid bodies, deformable objects, and fluids. *arXiv preprint arXiv:1810.01566*,
393 2018.
- 394 [36] X. Lin, C. Qi, Y. Zhang, Z. Huang, K. Fragkiadaki, Y. Li, C. Gan, and D. Held. Planning with
395 spatial-temporal abstraction from point clouds for deformable object manipulation. *arXiv*
396 *preprint arXiv:2210.15751*, 2022.
- 397 [37] Y. Wang, Y. Li, K. Driggs-Campbell, L. Fei-Fei, and J. Wu. Dynamic-resolution model
398 learning for object pile manipulation. *arXiv preprint arXiv:2306.16700*, 2023.
- 399 [38] H. Shi, H. Xu, S. Clarke, Y. Li, and J. Wu. Robocook: Long-horizon elasto-plastic object
400 manipulation with diverse tools. *arXiv preprint arXiv:2306.14447*, 2023.
- 401 [39] X. Lin, Y. Wang, Z. Huang, and D. Held. Learning visible connectivity dynamics for cloth
402 smoothing. In *Conference on Robot Learning*, pages 256–266. PMLR, 2022.
- 403 [40] J. Abou-Chakra, K. Rana, F. Dayoub, and N. Sünderhauf. Physically embodied gaussian
404 splatting: A realtime correctable world model for robotics. *arXiv preprint arXiv:2406.10788*,
405 2024.
- 406 [41] D. Bauer, Z. Xu, and S. Song. Doughnet: A visual predictive model for topological manipu-
407 lation of deformable objects. *arXiv preprint arXiv:2404.12524*, 2024.
- 408 [42] T. Schmidt, R. Newcombe, and D. Fox. Self-supervised visual descriptor learning for dense
409 correspondence. *IEEE Robotics and Automation Letters*, 2(2):420–427, 2016.
- 410 [43] P. R. Florence, L. Manuelli, and R. Tedrake. Dense object nets: Learning dense visual object
411 descriptors by and for robotic manipulation. *arXiv preprint arXiv:1806.08756*, 2018.
- 412 [44] T. D. Kulkarni, A. Gupta, C. Ionescu, S. Borgeaud, M. Reynolds, A. Zisserman, and V. Mnih.
413 Unsupervised learning of object keypoints for perception and control. *Advances in neural*
414 *information processing systems*, 32, 2019.
- 415 [45] Z. Qin, K. Fang, Y. Zhu, L. Fei-Fei, and S. Savarese. Keto: Learning keypoint representations
416 for tool manipulation. In *2020 IEEE International Conference on Robotics and Automation*
417 *(ICRA)*, pages 7278–7285. IEEE, 2020.
- 418 [46] P. Sundaresan, J. Grannen, B. Thananjeyan, A. Balakrishna, M. Laskey, K. Stone, J. E. Gon-
419 zalez, and K. Goldberg. Learning rope manipulation policies using dense object descriptors
420 trained on synthetic depth data. In *2020 IEEE International Conference on Robotics and*
421 *Automation (ICRA)*, pages 9411–9418. IEEE, 2020.
- 422 [47] L. Manuelli, Y. Li, P. Florence, and R. Tedrake. Keypoints into the future: Self-supervised
423 correspondence in model-based reinforcement learning. *arXiv preprint arXiv:2009.05085*,
424 2020.
- 425 [48] B. Chen, P. Abbeel, and D. Pathak. Unsupervised learning of visual 3d keypoints for control.
426 In *International Conference on Machine Learning*, pages 1539–1549. PMLR, 2021.

- 427 [49] A. Simeonov, Y. Du, Y.-C. Lin, A. R. Garcia, L. P. Kaelbling, T. Lozano-Pérez, and
428 P. Agrawal. Se (3)-equivariant relational rearrangement with neural descriptor fields. In
429 *Conference on Robot Learning*, pages 835–846. PMLR, 2023.
- 430 [50] M. Vecerik, C. Doersch, Y. Yang, T. Davchev, Y. Aytar, G. Zhou, R. Hadsell, L. Agapito, and
431 J. Scholz. Robotap: Tracking arbitrary points for few-shot visual imitation. *arXiv preprint*
432 *arXiv:2308.15975*, 2023.
- 433 [51] E. Chun, Y. Du, A. Simeonov, T. Lozano-Perez, and L. Kaelbling. Local neural descriptor
434 fields: Locally conditioned object representations for manipulation. In *2023 IEEE Interna-*
435 *tional Conference on Robotics and Automation (ICRA)*, pages 1830–1836. IEEE, 2023.
- 436 [52] S. Bahl, R. Mendonca, L. Chen, U. Jain, and D. Pathak. Affordances from human videos
437 as a versatile representation for robotics. In *Proceedings of the IEEE/CVF Conference on*
438 *Computer Vision and Pattern Recognition*, pages 13778–13790, 2023.
- 439 [53] C. Wen, X. Lin, J. So, K. Chen, Q. Dou, Y. Gao, and P. Abbeel. Any-point trajectory modeling
440 for policy learning. *arXiv preprint arXiv:2401.00025*, 2023.
- 441 [54] H. Bharadhwaj, R. Mottaghi, A. Gupta, and S. Tulsiani. Track2act: Predicting point tracks
442 from internet videos enables diverse zero-shot robot manipulation, 2024.
- 443 [55] Z. Kingston, M. Moll, and L. E. Kavraki. Sampling-based methods for motion planning with
444 constraints. *Annual review of control, robotics, and autonomous systems*, 1:159–185, 2018.
- 445 [56] N. Ratliff, M. Zucker, J. A. Bagnell, and S. Srinivasa. Chomp: Gradient optimization tech-
446 niques for efficient motion planning. In *2009 IEEE international conference on robotics and*
447 *automation*, pages 489–494. IEEE, 2009.
- 448 [57] J. Schulman, Y. Duan, J. Ho, A. Lee, I. Awwal, H. Bradlow, J. Pan, S. Patil, K. Goldberg,
449 and P. Abbeel. Motion planning with sequential convex optimization and convex collision
450 checking. *The International Journal of Robotics Research*, 33(9):1251–1270, 2014.
- 451 [58] B. Sundaralingam, S. K. S. Hari, A. Fishman, C. Garrett, K. Van Wyk, V. Blukis, A. Millane,
452 H. Oleynikova, A. Handa, F. Ramos, et al. Curobo: Parallelized collision-free robot motion
453 generation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*,
454 pages 8112–8119. IEEE, 2023.
- 455 [59] T. Marcucci, J. Umenberger, P. Parrilo, and R. Tedrake. Shortest paths in graphs of convex
456 sets. *SIAM Journal on Optimization*, 34(1):507–532, 2024.
- 457 [60] N. D. Ratliff, J. Issac, D. Kappler, S. Birchfield, and D. Fox. Riemannian motion policies.
458 *arXiv preprint arXiv:1801.02854*, 2018.
- 459 [61] M. Posa, S. Kuindersma, and R. Tedrake. Optimization and stabilization of trajectories for
460 constrained dynamical systems. In *2016 IEEE International Conference on Robotics and*
461 *Automation (ICRA)*, pages 1366–1373. IEEE, 2016.
- 462 [62] I. Mordatch, E. Todorov, and Z. Popović. Discovery of complex behaviors through contact-
463 invariant optimization. *ACM Transactions on Graphics (ToG)*, 31(4):1–8, 2012.
- 464 [63] I. Mordatch, Z. Popović, and E. Todorov. Contact-invariant optimization for hand manipula-
465 tion. In *Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer anima-*
466 *tion*, pages 137–144, 2012.
- 467 [64] M. Posa, C. Cantu, and R. Tedrake. A direct method for trajectory optimization of rigid
468 bodies through contact. *The International Journal of Robotics Research*, 33(1):69–81, 2014.

- 469 [65] T. Howell, N. Gileadi, S. Tunyasuvunakool, K. Zakka, T. Erez, and Y. Tassa. Predictive
470 sampling: Real-time behaviour synthesis with mujoco. *arXiv preprint arXiv:2212.00541*,
471 2022.
- 472 [66] Z. Liu, G. Zhou, J. He, T. Marcucci, F.-F. Li, J. Wu, and Y. Li. Model-based control with
473 sparse neural dynamics. *Advances in Neural Information Processing Systems*, 36, 2024.
- 474 [67] K. M. Lynch and M. T. Mason. Stable pushing: Mechanics, controllability, and planning. *The
475 international journal of robotics research*, 15(6):533–556, 1996.
- 476 [68] Y. Hou, Z. Jia, and M. T. Mason. Fast planning for 3d any-pose-reorienting using pivoting. In
477 *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1631–1638.
478 IEEE, 2018.
- 479 [69] J.-P. Sleiman, F. Farshidian, and M. Hutter. Versatile multicontact planning and control for
480 legged loco-manipulation. *Science Robotics*, 8(81):eadg5014, 2023.
- 481 [70] W. Yang and M. Posa. Dynamic on-palm manipulation via controlled sliding. *arXiv preprint
482 arXiv:2405.08731*, 2024.
- 483 [71] B. P. Graesdal, S. Y. Chia, T. Marcucci, S. Morozov, A. Amice, P. A. Parrilo, and
484 R. Tedrake. Towards tight convex relaxations for contact-rich manipulation. *arXiv preprint
485 arXiv:2402.10312*, 2024.
- 486 [72] K. Yunt and C. Glocker. Trajectory optimization of mechanical hybrid systems using sumt.
487 In *9th IEEE International Workshop on Advanced Motion Control, 2006.*, pages 665–671.
488 IEEE, 2006.
- 489 [73] K. Yunt and C. Glocker. A combined continuation and penalty method for the determination
490 of optimal hybrid mechanical trajectories. In *Iutam Symposium on Dynamics and Control of
491 Nonlinear Systems with Uncertainty: Proceedings of the IUTAM Symposium held in Nanjing,
492 China, September 18-22, 2006*, pages 187–196. Springer, 2007.
- 493 [74] K. Yunt. An augmented lagrangian based shooting method for the optimal trajectory gen-
494 eration of switching lagrangian systems. *Dynamics of Continuous, Discrete and Impulsive
495 Systems Series B: Applications and Algorithms*, 18(5):615–645, 2011.
- 496 [75] F. Lagriffoul, D. Dimitrov, A. Saffiotti, and L. Karlsson. Constraint propagation on interval
497 bounds for dealing with geometric backtracking. In *2012 IEEE/RSJ International Conference
498 on Intelligent Robots and Systems*, pages 957–964. IEEE, 2012.
- 499 [76] F. Lagriffoul, D. Dimitrov, J. Bidot, A. Saffiotti, and L. Karlsson. Efficiently combining
500 task and motion planning using geometric constraints. *The International Journal of Robotics
501 Research*, 33(14):1726–1747, 2014.
- 502 [77] T. Lozano-Pérez and L. P. Kaelbling. A constraint-based method for solving sequential ma-
503 nipulation planning problems. In *2014 IEEE/RSJ International Conference on Intelligent
504 Robots and Systems*, pages 3684–3691. IEEE, 2014.
- 505 [78] Z. Yang, J. Mao, Y. Du, J. Wu, J. B. Tenenbaum, T. Lozano-Pérez, and L. P. Kaelbling. Com-
506 positional diffusion-based continuous constraint solvers. *arXiv preprint arXiv:2309.00966*,
507 2023.
- 508 [79] T. Silver, R. Chitnis, J. Tenenbaum, L. P. Kaelbling, and T. Lozano-Pérez. Learning sym-
509 bolic operators for task and motion planning. In *2021 IEEE/RSJ International Conference on
510 Intelligent Robots and Systems (IROS)*, pages 3182–3189. IEEE, 2021.
- 511 [80] B. Vu, T. Migimatsu, and J. Bohg. Coast: Constraints and streams for task and motion
512 planning. *arXiv preprint arXiv:2405.08572*, 2024.

- 513 [81] M. Toussaint. Logic-geometric programming: An optimization-based approach to combined
514 task and motion planning. In *Twenty-Fourth International Joint Conference on Artificial
515 Intelligence*, 2015.
- 516 [82] M. Toussaint and M. Lopes. Multi-bound tree search for logic-geometric programming in
517 cooperative manipulation domains. In *2017 IEEE International Conference on Robotics and
518 Automation (ICRA)*, pages 4044–4051. IEEE, 2017.
- 519 [83] M. A. Toussaint, K. R. Allen, K. A. Smith, and J. B. Tenenbaum. Differentiable physics
520 and stable modes for tool-use and manipulation planning. *Robotics: Science and Systems
521 Foundation*, 2018.
- 522 [84] J.-S. Ha, D. Driess, and M. Toussaint. A probabilistic framework for constrained manipula-
523 tions and task and motion planning under uncertainty. In *2020 IEEE International Conference
524 on Robotics and Automation (ICRA)*, pages 6745–6751. IEEE, 2020.
- 525 [85] T. Xue, A. Razmjoo, and S. Calinon. D-lgp: Dynamic logicgeometric program for reactive
526 task and motion planning. *arXiv preprint arXiv:2312.02731*, 2023.
- 527 [86] D. Driess, O. Oguz, J.-S. Ha, and M. Toussaint. Deep visual heuristics: Learning feasibility
528 of mixed-integer programs for manipulation planning. In *2020 IEEE international conference
529 on robotics and automation (ICRA)*, pages 9563–9569. IEEE, 2020.
- 530 [87] G. Sutanto, I. R. Fernández, P. Englert, R. K. Ramachandran, and G. Sukhatme. Learning
531 equality constraints for motion planning on manifolds. In *Conference on Robot Learning*,
532 pages 2292–2305. PMLR, 2021.
- 533 [88] Z. Yang, C. R. Garrett, T. Lozano-Pérez, L. Kaelbling, and D. Fox. Sequence-based plan
534 feasibility prediction for efficient task and motion planning. *arXiv preprint arXiv:2211.01576*,
535 2022.
- 536 [89] G. S. Camps, R. Dyro, M. Pavone, and M. Schwager. Learning deep sdf maps online for
537 robot navigation and exploration. *arXiv preprint arXiv:2207.10782*, 2022.
- 538 [90] Y. Hu, Q. Xie, V. Jain, J. Francis, J. Patrikar, N. Keetha, S. Kim, Y. Xie, T. Zhang, Z. Zhao,
539 et al. Toward general-purpose robots via foundation models: A survey and meta-analysis.
540 *arXiv preprint arXiv:2312.08782*, 2023.
- 541 [91] R. Firoozi, J. Tucker, S. Tian, A. Majumdar, J. Sun, W. Liu, Y. Zhu, S. Song, A. Kapoor,
542 K. Hausman, et al. Foundation models in robotics: Applications, challenges, and the future.
543 *arXiv preprint arXiv:2312.07843*, 2023.
- 544 [92] K. Kawaharazuka, T. Matsushima, A. Gambardella, J. Guo, C. Paxton, and A. Zeng. Real-
545 world robot applications of foundation models: A review. *arXiv preprint arXiv:2402.05741*,
546 2024.
- 547 [93] S. Yang, O. Nachum, Y. Du, J. Wei, P. Abbeel, and D. Schuurmans. Foundation models for
548 decision making: Problems, methods, and opportunities. *arXiv preprint arXiv:2303.04129*,
549 2023.
- 550 [94] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
551 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language super-
552 vision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- 553 [95] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever.
554 Zero-shot text-to-image generation. In *International conference on machine learning*, pages
555 8821–8831. Pmlr, 2021.

- 556 [96] J. Li, D. Li, C. Xiong, and S. Hoi. Blip: Bootstrapping language-image pre-training for uni-
557 fied vision-language understanding and generation. In *International conference on machine*
558 *learning*, pages 12888–12900. PMLR, 2022.
- 559 [97] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida,
560 J. Altschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint*
561 *arXiv:2303.08774*, 2023.
- 562 [98] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with
563 frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- 564 [99] H. Huang, F. Lin, Y. Hu, S. Wang, and Y. Gao. Copa: General robotic manipulation through
565 spatial constraints of parts with foundation models. *arXiv preprint arXiv:2403.08248*, 2024.
- 566 [100] F. Liu, K. Fang, P. Abbeel, and S. Levine. Moka: Open-vocabulary robotic manipulation
567 through mark-based visual prompting. *arXiv preprint arXiv:2403.03174*, 2024.
- 568 [101] S. Nasiriany, F. Xia, W. Yu, T. Xiao, J. Liang, I. Dasgupta, A. Xie, D. Driess, A. Wahid,
569 Z. Xu, et al. Pivot: Iterative visual prompting elicits actionable knowledge for vlms. *arXiv*
570 *preprint arXiv:2402.07872*, 2024.
- 571 [102] Y. Hu, F. Lin, T. Zhang, L. Yi, and Y. Gao. Look before you leap: Unveiling the power of
572 gpt-4v in robotic vision-language planning. *arXiv preprint arXiv:2311.17842*, 2023.
- 573 [103] Y. Du, M. Yang, P. Florence, F. Xia, A. Wahid, B. Ichter, P. Sermanet, T. Yu, P. Abbeel, J. B.
574 Tenenbaum, et al. Video language planning. *arXiv preprint arXiv:2310.10625*, 2023.
- 575 [104] Y. Hong, H. Zhen, P. Chen, S. Zheng, Y. Du, Z. Chen, and C. Gan. 3d-llm: Injecting the 3d
576 world into large language models. *Advances in Neural Information Processing Systems*, 36:
577 20482–20494, 2023.
- 578 [105] B. Chen, Z. Xu, S. Kirmani, B. Ichter, D. Sadigh, L. Guibas, and F. Xia. Spatialvlm: En-
579 dowing vision-language models with spatial reasoning capabilities. In *Proceedings of the*
580 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465,
581 2024.
- 582 [106] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei. Voxposer: Composable 3d
583 value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*,
584 2023.
- 585 [107] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess,
586 A. Dubey, C. Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to
587 robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- 588 [108] J. Gao, B. Sarkar, F. Xia, T. Xiao, J. Wu, B. Ichter, A. Majumdar, and D. Sadigh.
589 Physically grounded vision-language models for robotic manipulation. *arXiv preprint*
590 *arXiv:2309.02561*, 2023.
- 591 [109] Y. Wang, T.-H. Wang, J. Mao, M. Hagenow, and J. Shah. Grounding language plans in demon-
592 strations through counterfactual perturbations. *arXiv preprint arXiv:2403.17124*, 2024.
- 593 [110] J. Hsu, J. Mao, and J. Wu. Ns3d: Neuro-symbolic grounding of 3d objects and relations.
594 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
595 pages 2614–2623, 2023.
- 596 [111] J. Gao, B. Sarkar, F. Xia, T. Xiao, J. Wu, B. Ichter, A. Majumdar, and D. Sadigh. Physi-
597 cally grounded vision-language models for robotic manipulation. In *2024 IEEE International*
598 *Conference on Robotics and Automation (ICRA)*, pages 12462–12469. IEEE, 2024.

- 599 [112] W. Yuan, J. Duan, V. Blukis, W. Pumacay, R. Krishna, A. Murali, A. Mousavian, and D. Fox.
600 Robopoint: A vision-language model for spatial affordance prediction for robotics. *arXiv*
601 *preprint arXiv:2406.10721*, 2024.
- 602 [113] J. Duan, W. Yuan, W. Pumacay, Y. R. Wang, K. Ehsani, D. Fox, and R. Krishna. Manipulate-
603 anything: Automating real-world robots using vision-language models. *arXiv preprint*
604 *arXiv:2406.18915*, 2024.
- 605 [114] S. Tong, Z. Liu, Y. Zhai, Y. Ma, Y. LeCun, and S. Xie. Eyes wide shut? exploring the visual
606 shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer*
607 *Vision and Pattern Recognition*, pages 9568–9578, 2024.
- 608 [115] T. Thrush, R. Jiang, M. Bartolo, A. Singh, A. Williams, D. Kiela, and C. Ross. Winoground:
609 Probing vision and language models for visio-linguistic compositionality. In *Proceedings of*
610 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248,
611 2022.
- 612 [116] M. Yuksekogonul, F. Bianchi, P. Kalluri, D. Jurafsky, and J. Zou. When and why vision-
613 language models behave like bags-of-words, and what to do about it? In *The Eleventh Inter-*
614 *national Conference on Learning Representations*, 2023.
- 615 [117] C.-Y. Hsieh, J. Zhang, Z. Ma, A. Kembhavi, and R. Krishna. Sugarcrepe: Fixing hackable
616 benchmarks for vision-language compositionality. *Advances in neural information processing*
617 *systems*, 36, 2024.
- 618 [118] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerg-
619 ing properties in self-supervised vision transformers. In *Proceedings of the International*
620 *Conference on Computer Vision (ICCV)*, 2021.
- 621 [119] S. Amir, Y. Gandelsman, S. Bagon, and T. Dekel. Deep vit features as dense visual descrip-
622 tors. *arXiv preprint arXiv:2112.05814*, 2(3):4, 2021.
- 623 [120] L. Melas-Kyriazi, C. Rupprecht, I. Laina, and A. Vedaldi. Deep spectral methods: A sur-
624 prisingly strong baseline for unsupervised semantic segmentation and localization. In *Pro-*
625 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
626 8364–8375, 2022.
- 627 [121] Y. Wang, Z. Li, M. Zhang, K. Driggs-Campbell, J. Wu, L. Fei-Fei, and Y. Li. D³fields:
628 Dynamic 3d descriptor fields for zero-shot generalizable robotic manipulation. *arXiv preprint*
629 *arXiv:2309.16118*, 2023.
- 630 [122] X. Lin, J. So, S. Mahalingam, F. Liu, and P. Abbeel. Spawnnet: Learning generalizable
631 visuomotor skills from pre-trained networks. *arXiv preprint arXiv:2307.03567*, 2023.
- 632 [123] N. Di Palo and E. Johns. Keypoint action tokens enable in-context imitation learning in
633 robotics. *arXiv preprint arXiv:2403.19578*, 2024.
- 634 [124] N. Di Palo and E. Johns. Dinobot: Robot manipulation via retrieval and alignment with vision
635 foundation models. *arXiv preprint arXiv:2402.13181*, 2024.
- 636 [125] O. Y. Lee, A. Xie, K. Fang, K. Pertsch, and C. Finn. Affordance-guided reinforcement learn-
637 ing via visual prompting. *arXiv preprint arXiv:2407.10341*, 2024.
- 638 [126] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau,
639 E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk,
640 M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard,
641 T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with
642 NumPy. *Nature*, 585(7825):357–362, Sept. 2020. doi:10.1038/s41586-020-2649-2. URL
643 <https://doi.org/10.1038/s41586-020-2649-2>.

- 644 [127] R. Tedrake. *Underactuated Robotics*. 2023. URL <https://underactuated.csail.mit.edu>.
645
- 646 [128] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau,
647 E. Burovski, P. Peterson, W. Weckesser, J. Bright, et al. Scipy 1.0: fundamental algorithms
648 for scientific computing in python. *Nature methods*, 17(3):261–272, 2020.
- 649 [129] Y. Xiang, D. Sun, W. Fan, and X. Gong. Generalized simulated annealing algorithm and its
650 application to the thomson model. *Physics Letters A*, 233(3):216–220, 1997.
- 651 [130] D. Kraft. A software package for sequential quadratic programming. *Forschungsbericht-*
652 *Deutsche Forschungs- und Versuchsanstalt für Luft- und Raumfahrt*, 1988.
- 653 [131] H.-S. Fang, C. Wang, H. Fang, M. Gou, J. Liu, H. Yan, W. Liu, Y. Xie, and C. Lu. Anygrasp:
654 Robust and efficient grasp perception in spatial and temporal domains. *IEEE Transactions on*
655 *Robotics*, 2023.
- 656 [132] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead,
657 A. C. Berg, W.-Y. Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- 658 [133] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation for games, robotics
659 and machine learning. 2016.
- 660 [134] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. Viola: Imitation learning for vision-based manipulation
661 with object proposal priors. *6th Annual Conference on Robot Learning*, 2022.
- 662 [135] M. Minderer, A. Gritsenko, A. Stone, M. Neumann, D. Weissenborn, A. Dosovitskiy, A. Ma-
663 hendran, A. Arnab, M. Dehghani, Z. Shen, et al. Simple open-vocabulary object detection
664 with vision transformers. *arXiv preprint arXiv:2205.06230*, 2022.
- 665 [136] H. K. Cheng, S. W. Oh, B. Price, J.-Y. Lee, and A. Schwing. Putting the object back into
666 video object segmentation. In *arXiv*, 2023.
- 667 [137] T. Darcet, M. Oquab, J. Mairal, and P. Bojanowski. Vision transformers need registers. *arXiv*
668 *preprint arXiv:2309.16588*, 2023.
- 669 [138] D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis.
670 *IEEE Transactions on pattern analysis and machine intelligence*, 24(5):603–619, 2002.
- 671 [139] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
672 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
673 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
674 *Learning Research*, 12:2825–2830, 2011.
- 675 [140] J. Yang, H. Zhang, F. Li, X. Zou, C. Li, and J. Gao. Set-of-mark prompting unleashes ex-
676 traordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*, 2023.
- 677 [141] A. W. Harley, Z. Fang, and K. Fragkiadaki. Particle video revisited: Tracking through oc-
678 clusions using point trajectories. In *European Conference on Computer Vision*, pages 59–75.
679 Springer, 2022.
- 680 [142] Q. Wang, Y.-Y. Chang, R. Cai, Z. Li, B. Hariharan, A. Holynski, and N. Snavely. Tracking
681 everything everywhere all at once. In *Proceedings of the IEEE/CVF International Conference*
682 *on Computer Vision*, pages 19795–19806, 2023.
- 683 [143] Y. Zheng, A. W. Harley, B. Shen, G. Wetzstein, and L. J. Guibas. Pointodyssey: A large-scale
684 synthetic dataset for long-term point tracking. In *Proceedings of the IEEE/CVF International*
685 *Conference on Computer Vision*, pages 19855–19865, 2023.

- 686 [144] N. Karaev, I. Rocco, B. Graham, N. Neverova, A. Vedaldi, and C. Rupprecht. Cotracker: It is
687 better to track together. *arXiv preprint arXiv:2307.07635*, 2023.
- 688 [145] C. Doersch, Y. Yang, M. Vecerik, D. Gokay, A. Gupta, Y. Aytar, J. Carreira, and A. Zisserman.
689 Tapir: Tracking any point with per-frame initialization and temporal refinement. In *Proceed-*
690 *ings of the IEEE/CVF International Conference on Computer Vision*, pages 10061–10072,
691 2023.
- 692 [146] C. Doersch, Y. Yang, D. Gokay, P. Luc, S. Koppula, A. Gupta, J. Heyward, R. Goroshin,
693 J. Carreira, and A. Zisserman. Bootstap: Bootstrapped training for tracking-any-point. *arXiv*
694 *preprint arXiv:2402.00847*, 2024.
- 695 [147] Y. Xiao, Q. Wang, S. Zhang, N. Xue, S. Peng, Y. Shen, and X. Zhou. Spatialtracker: Tracking
696 any 2d pixels in 3d space. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
697 *and Pattern Recognition*, pages 20406–20417, 2024.
- 698 [148] J. Luiten, G. Kopanas, B. Leibe, and D. Ramanan. Dynamic 3d gaussians: Tracking by
699 persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713*, 2023.
- 700 [149] A. Millane, H. Oleynikova, E. Wirbel, R. Steiner, V. Ramasamy, D. Tingdahl, and R. Sieg-
701 wart. nvblox: Gpu-accelerated incremental signed distance field mapping. *arXiv preprint*
702 *arXiv:2311.00626*, 2023.
- 703 [150] X. Li, M. Zhang, Y. Geng, H. Geng, Y. Long, Y. Shen, R. Zhang, J. Liu, and H. Dong. Ma-
704 nipllm: Embodied multimodal large language model for object-centric robotic manipulation.
705 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
706 pages 18061–18070, 2024.
- 707 [151] W. Xia, D. Wang, X. Pang, Z. Wang, B. Zhao, and D. Hu. Kinematic-aware prompting for
708 generalizable articulated object manipulation with llms. *arXiv preprint arXiv:2311.02847*,
709 2023.
- 710 [152] S. Huang, H. Chang, Y. Liu, Y. Zhu, H. Dong, P. Gao, A. Boularias, and H. Li. A3vlm: Ac-
711 tionable articulation-aware vision language model. *arXiv preprint arXiv:2406.07549*, 2024.
- 712 [153] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine,
713 M. Lingelbach, J. Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday
714 activities and realistic simulation. In *Conference on Robot Learning*, pages 80–93. PMLR,
715 2023.
- 716 [154] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
717 I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*,
718 30, 2017.
- 719 [155] A. Goyal, J. Xu, Y. Guo, V. Blukis, Y.-W. Chao, and D. Fox. Rvt: Robotic view transformer
720 for 3d object manipulation. In *Conference on Robot Learning*, pages 694–710. PMLR, 2023.
- 721 [156] A. Goyal, V. Blukis, J. Xu, Y. Guo, Y.-W. Chao, and D. Fox. Rvt-2: Learning precise manip-
722 ulation from few demonstrations. *arXiv preprint arXiv:2406.08545*, 2024.
- 723 [157] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. De-
724 hghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Trans-
725 formers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

726 **A Appendix**

727 **A.1 Pseudo-code for Sequential Manipulation with Relational Keypoint Constraints**

Algorithm 1 Relational Keypoint Constraints for Sequential Manipulation

```

1: Initialize current stage  $i \leftarrow 1$ , and current time  $t \leftarrow 1$ 
2: while  $i \leq N$  do
3:   if  $\exists f \in \mathcal{C}_{\text{path}}^{(i)}$  s.t.  $f(k_t) > 0$  then
4:      $i \leftarrow i - 1$ 
5:     continue
6:   end if
7:   if  $\text{distance}(\mathbf{e}_t, \mathbf{e}_{g_i}) < \epsilon$  then
8:      $i \leftarrow i + 1$ 
9:     continue
10:  end if
11:  Solve sub-goal problem for stage  $i$  to obtain  $\mathbf{e}_{g_i}$  (Eq. 2)
12:  Solve path problem for stage  $i$  to obtain  $\mathbf{e}_{t:g_i}$  (Eq. 3)
13:  Execute the next  $m$  actions  $\mathbf{e}_{t+1:t+m}$ 
14:   $t \leftarrow t + m + 1$ 
15: end while

```

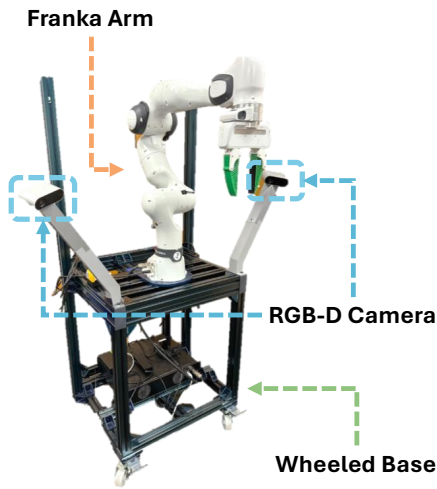


Figure 6: Wheeled Single-Arm Platform.

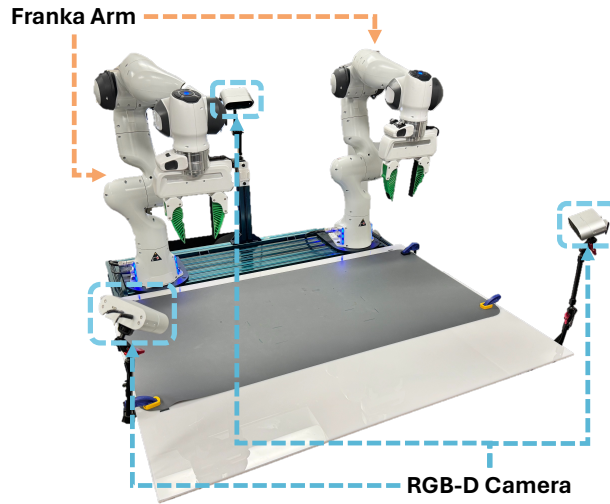


Figure 7: Stationary Dual-Arm Platform.

728 **A.2 Wheeled Single-Arm Platform**

729 One of our investigated platform is a Franka arm mounted on a wheeled base built with Vention
730 frames (shown in Figure 6). Note that the base does not have motors and thus cannot move au-
731 tonomously, but its mobility nevertheless allows us to investigate the proposed method outside of
732 lab environments.

733 Since our pipeline produces a sequence of 6-DoF end-effector poses, we use position control in all
734 experiments, which is running at a fixed frequency of 20 Hz. Specifically, once the robot is given a
735 target end-effector pose in the world frame, we first clip the pose to the pre-defined workspace. Then
736 we linearly interpolate from the current pose to the target pose with a step size of 5mm for position
737 and 1 degree for rotation. To move to each interpolated pose, we first calculate inverse kinematics
738 to obtain the target joint positions based on current joint positions (IK solver from PyBullet [133]).
739 Then we use the joint impedance controller from Deoxys [134] to reach to the target joint positions.

740 Two RGB-D cameras, Orbbec Femto Bolt, are mounted on each side of the robot facing the
741 workspace center. The cameras capture RGB images and point clouds at a fixed frequency of 20
742 Hz.

743 A.3 Stationary Dual-Arm Platform

744 We also investigate the method on a stationary dual-arm platform consisting of two Franka arms
745 mounted in front of a tabletop workspace (shown in Figure 7). We share the same controller as the
746 wheeled single-arm platform with the exception that the two arms are controlled simultaneously at
747 20 Hz. Specifically, our pipeline jointly solves two 6-DoF end-effector pose sequences, which are
748 sent to the low-level controller together. The controller subsequently calculates IK for both arms
749 and moves the arms using joint impedance control.

750 Three RGB-D cameras, Orbbec Femto Bolt, are mounted on this platform. Two cameras are
751 mounted on the left and right sides and one camera is mounted in the back. The cameras simi-
752 larly capture RGB images and point clouds at a fixed frequency of 20 Hz.

753 A.4 Evaluation Details

754 Below we discuss the evaluation details for the experiments reported in Section 4.1 and Section 4.2.

755 A.4.1 Details for In-the-Wild and Bimanual Manipulation (Section 4.1)

756 For each task, 10 initial different configurations of objects are selected, which cover the full
757 workspace but are manually verified to ensure they are kinematically feasible for the robot. For
758 each trial, a human operator restores the scene to the corresponding configuration and initiates the
759 system. Due to the challenge of developing automatic success criteria for the diverse set of objects
760 and environments investigated in this work, success rates are measured by the operator with the cri-
761 terion reported under each task description below. For experiments involving external disturbances,
762 the set of disturbances for all trials is pre-selected, and one disturbance is applied to each trial.
763 Specifically, the disturbance is introduced by a human operator using hands to change the object’s
764 pose. Collision checking is disabled for all tasks involving deformable objects.

765 **Pour Tea:** The environment consists of a teapot and a cup placed on a counter table in a kitchen
766 setting. The task involves three stages: grasping the handle, aligning the teapot to the top of the cup,
767 and pouring the tea into the cup. The success criterion requires that the teapot remains upright until
768 the pouring stage, and at the end, the spout must be aligned and tilted on top of the cup opening.

769 **Recycle Can:** The environment includes one of three types of cans (Coke, Zero Coke, Zero Sprite),
770 a recycle bin with a narrow opening (such that the cans may only go in when they are upright), a
771 landfill bin, and a compost bin, all situated inside an office building. The task involves two stages:
772 grasping the can and reorienting it on top of the recycle bin before dropping it. The success criterion
773 is that the can is successfully thrown into the bin.

774 **Stow Book:** The environment consists of a target book placed on a side table and a real-size book-
775 shelf with a 15cm opening among the placed books, all inside an office environment. The task
776 involves two stages: grasping the target book on the side and stowing it inside the opening in the
777 shelf. The success criterion is that the target book is placed steadily after the robot releases the
778 gripper, and the robot must not bump into the shelf or other placed books.

779 **Tape Box:** The environment includes a cardboard box, a packaging tape with a dispenser sitting on
780 top of the box that already has one side taped, and a human user collaborating with the robot. The
781 tape has already been unrolled to be enough for taping because unrolling typically requires a large
782 force that exceeds the limit of the robot arm. The task involves two stages: while a human operator
783 is squeezing the box, the robot needs to grasp the tape and align it to the correct side to complete
784 the taping. The success criterion is that the tape must end up in the correct position such that it is
785 aligned with the seam.

786 **(Bimanual) Fold Garment:** The environment consists of a sweater placed flat close to the
787 workspace center, with small deformations on the sleeves, neck, and bottom. The task typically
788 requires four stages: grasping both sleeves, folding them to the middle, grasping the neck, and fold-
789 ing it to the bottom. The success criterion does not enforce consistent stages; as long as the sweater
790 is folded such that it occupies at most half of the original surface size, it is regarded as a success.

791 **(Bimanual) Pack Shoes:** The environment includes an empty shoe box placed close to the
792 workspace center, with two shoes placed on opposite sides of the box in random poses. The task
793 involves two stages: grasping the shoes simultaneously and placing them in the shoe box. The suc-
794 cess criterion does not enforce consistent stages; as long as the shoes are placed into the box without
795 being stacked together or causing bimanual self-collision, it is considered successful.

796 **(Bimanual) Collaborative Folding:** The environment consists of a large blanket (pre-folded to an
797 appropriate size that occupies about 70% of the workspace due to its size exceeding the workspace
798 limit) and a human user collaborating with the robot. The task involves two stages: the robot must
799 grasp the two corners of the blanket opposite to the human user, and the second stage is aligning
800 the two corners with the two corners that the human has grasped. The success criterion is that the
801 robot has grasped the correct corners and can align them with the correct human arms (left-left,
802 right-right).

803 **A.4.2 Details on Baseline Methods**

804 We use VoxPoser [106] as the main baseline method as it makes similar assumptions that no task-
805 specific data or pre-defined motion primitives are required. We adapt VoxPoser to our setting with
806 certain modifications to ensure fairness in comparisons. Specifically, we use the same VLM, GPT-
807 4o [6], that takes the same camera input. We also augment the original prompt from the paper with
808 the prompt used in this work to ensure it has sufficient context. We only use the affordance, rotation,
809 and gripper action value maps and ignore the avoidance and velocity value maps because they are
810 not necessary for our tasks. We also only consider the scenario where the “entity of interest” is
811 the robot end-effector instead of objects in the scene. The latter is tailored for pushing task, which
812 is not being studied in this work. We use OWL-ViT [135] for open-vocabulary object detection,
813 SAM [132] for initial-frame segmentation, and Cutie [136] for mask tracking.

814 **A.4.3 Details for Generalization in Manipulation Strategies (Section 4.2)**

815 The dual-arm robot is tasked with folding eight different categories of clothing. We use two metrics
816 for evaluation: “Strategy Success” and “Execution Success,” where the former evaluates whether
817 keypoints are proposed and constraints are written appropriately, and the latter evaluates the robotic
818 system’s execution given successful strategies.

819 To evaluate “Strategy Success,” the garment is initialized close to the center of the workspace. A
820 back-mounted RGB-D camera captures the RGB image. Then, the keypoint proposal module gener-
821 ates keypoint candidates using the captured image, which are then overlaid on top of the original
822 image with numerical marks $\{0, \dots, K - 1\}$. The overlaid image, along with the same generic
823 prompt, is fed into GPT-4 [6] to generate the ReKep constraints. Since folding garments is itself an
824 open-ended problem without ground-truth strategies, we manually judge if the proposed keypoints
825 and the generated constraints are correct. Note that since the constraints are to be executed by a
826 bimanual robot, and the constraints are almost always connecting (folding) two keypoints such that
827 they are aligned, correctness is measured by whether it is (potentially) executable by the robot with-
828 out causing self-collision (arms crossing over to opposite sides) and whether the folding strategy can
829 fold the garment to at most half of its original surface area.

830 To evaluate “Execution Success,” we take the generated strategies in the previous section that are
831 marked as successful for each garment and execute the sequence on the dual-arm platform, with a
832 total of 10 trials for each garment. Point tracking is disabled as we observe that our point tracker
833 predicts unstable tracks when the garment is potentially folded many times. Success is measured by
834 whether the garment is folded such that its surface area is at most half of its original surface area.

835 A.5 Implementation Details of Keypoint Proposal

836 Herein we describe how keypoint candidates in a scene are generated. For each platform, we use
837 one of the mounted RGB-D cameras to capture an image of size $h \times w \times 3$, depending on which
838 camera has the best holistic view of the environment, as all the keypoints need to be present in the
839 first frame for the proposed method. Given the captured image, we first use DINOv2 with registers
840 (ViT-S14) [5, 137] to extract the patch-wise features $\mathbf{F}_{\text{patch}} \in \mathbb{R}^{h' \times w' \times d}$. Then we perform bilinear
841 interpolation to upsample the features to the original image size, $\mathbf{F}_{\text{interp}} \in \mathbb{R}^{h \times w \times d}$. To ensure the
842 proposal covers all relevant objects in the scene, we extract all masks $\mathbf{M} = \{\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_n\}$ in
843 the scene using Segment Anything (SAM) [132]. Within each mask \mathbf{m}_i , we apply PCA to project the
844 features to three dimensions, $\mathbf{F}_{\text{PCA}} = \text{PCA}(\mathbf{F}_{\text{resized}}[\mathbf{m}_i], 3)$. We find that applying PCA improves
845 the clustering as it often removes details and artifacts related to texture that are not useful for our
846 tasks. For each mask j , we cluster the masked features $\mathbf{F}_{\text{interp}}[\mathbf{m}_j]$ using k -means with $k = 5$ with
847 the Euclidean distance metric. The median centroids of the clusters are used as keypoint candidates,
848 which are projected to a world coordinate \mathbb{R}^3 using a calibrated RGB-D camera. Note that we
849 also store which keypoint candidates originate from the same mask, which is later used as part of
850 the rigidity assumption in the optimization loops described in Sec. 3.3. Candidates outside of the
851 workspace bounds are filtered out. To avoid many points cluttered in a small region, we additionally
852 use Mean Shift [138, 139] (with a bandwidth 8cm) to filter out points that are close to each other.
853 Finally, the centroids are taken as final candidates. Alternatively, one may develop a pipeline using
854 only segmentation models [132, 140], but we leave comparisons to future work.

855 A.6 Querying Vision-Language Model

856 After we obtain the keypoint candidates, they are overlaid on the captured RGB image with numer-
857 ical marks $\{0, \dots, K - 1\}$. Then the image and the task instruction are fed into a vision-language
858 model with the prompt described below. The prompt contains only generic instructions with no
859 image-text in-context examples, although a few text-based examples are given to concretely ex-
860 plain the proposed method and the expected output from the model. Note that the majority of the
861 investigated tasks are not discussed in the provided prompt. As a result, the VLM is tasked with
862 generating ReKep constraints by leveraging its internalized world knowledge.

863 For the experiments conducted in this work, we use GPT-4o [6] as it is one of the latest available
864 models at the time of the experiments. However, due to rapid advancement in this field, the pipeline
865 can directly benefit from newer models that have better vision-language reasoning. Correspond-
866 ingly, we observe different models exhibit different behaviors when given the same prompt (with
867 the observation that newer models typically require less fine-grained instructions). As a result, in-
868 stead of developing the best prompt for the suite of tasks in this work, we focus on demonstrating
869 a full-stack pipeline consisting a key component that can be automated and continuously improved
870 by future development.

```
871 ## Instructions
872 Suppose you are controlling a robot to perform manipulation tasks by writing constraint functions in Python.
873 The manipulation task is given as an image of the environment, overlaid with keypoints marked with
874 their indices, along with a text instruction. The instruction starts with a parenthesis indicating
875 whether the robot has a single arm or is bimanual. For each given task, please perform the following
876 steps:
877 - Determine how many stages are involved in the task. Grasping must be an independent stage. Some examples:
878 - "(single-arm) pouring tea from teapot":
879   - 3 stages: "grasp teapot", "align teapot with cup opening", and "pour liquid"
880 - "(single-arm) put red block on top of blue block":
881   - 3 stages: "grasp red block", "align red block on top of blue block", and "release red block"
882 - "(bimanual) fold sleeves to the center":
883   - 2 stages: "left arm grasps left sleeve and right arm grasps right sleeve" and "both arms fold sleeves to
884     the center"
885 - "(bimanual) fold a jacket":
886   - 3 stages: "left arm grasps left sleeve and right arm grasps right sleeve", "both arms fold sleeves to
887     the center", and "grasp the neck with one arm (the other arm stays in place)", and "align the neck
888     with the bottom"
889 - For each stage, write two kinds of constraints, "sub-goal constraints" and "path constraints". The "sub-goal
890   constraints" are constraints that must be satisfied **at the end of the stage**, while the "path
891   constraints" are constraints that must be satisfied **within the stage**. Some examples:
892 - "(single-arm) pouring liquid from teapot":
893   - "grasp teapot" stage:
894     - sub-goal constraints: "align the end-effector with the teapot handle"
895     - path constraints: None
```

```

897 - "align teapot with cup opening" stage:
898 - sub-goal constraints: "the teapot spout needs to be 10cm above the cup opening"
899 - path constraints: "robot is grasping the teapot", and "the teapot must stay upright to avoid spilling"
900 - "pour liquid" stage:
901 - sub-goal constraints: "the teapot spout needs to be 5cm above the cup opening", "the teapot spout must
902   be tilted to pour liquid"
903 - path constraints: "the teapot spout is directly above the cup opening"
904 - (bimanual) fold sleeves to the center":
905 - "left arm grasps left sleeve and right arm grasps right sleeve" stage:
906 - sub-goal constraints: "left arm grasps left sleeve", "right arm grasps right sleeve"
907 - path constraints: None
908 - "both arms fold sleeves to the center" stage:
909 - sub-goal constraints: "left sleeve aligns with the center", "right sleeve aligns with the center"
910 - path constraints: None
911
912 Note:
913 - Each constraint takes a dummy end-effector point and a set of keypoints as input and returns a numerical
914   cost, where the constraint is satisfied if the cost is smaller than or equal to zero.
915 - For each stage, you may write 0 or more sub-goal constraints and 0 or more path constraints.
916 - Avoid using "if" statements in your constraints.
917 - Avoid using path constraints when manipulating deformable objects (e.g., clothing, towels).
918 - You do not need to consider collision avoidance. Focus on what is necessary to complete the task.
919 - Inputs to the constraints are as follows:
920 - `end_effector`: np.array of shape `(3,)` representing the end-effector position.
921 - `keypoints`: np.array of shape `(K, 3)` representing the keypoint positions.
922 - Inside of each function, you may use native Python functions and NumPy functions.
923 - For grasping stage, you should only write one sub-goal constraint that associates the end-effector with a
924   keypoint. No path constraints are needed.
925 - For non-grasping stage, you should not refer to the end-effector position.
926 - In order to move a keypoint, its associated object must be grasped in one of the previous stages.
927 - The robot can only grasp one object at a time.
928 - Grasping must be an independent stage from other stages.
929 - You may use two keypoints to form a vector, which can be used to specify a rotation (by specifying the angle
930   between the vector and a fixed axis).
931 - You may use multiple keypoints to specify a surface or volume.
932 - You may also use the center of multiple keypoints to specify a position.
933 - A single folding action should consist of two stages: one grasp and one place.
934
935 Structure your output in a single python code block as follows for single-arm robot:
936 ```python
937
938 # Your explanation of how many stages are involved in the task and what each stage is about.
939 # ...
940
941 num_stages = ?
942
943 ### stage 1 sub-goal constraints (if any)
944 def stage1_subgoal_constraint1(end_effector, keypoints):
945     """Put your explanation here."""
946     ...
947     return cost
948 # Add more sub-goal constraints if needed
949
950 ### stage 1 path constraints (if any)
951 def stage1_path_constraint1(end_effector, keypoints):
952     """Put your explanation here."""
953     ...
954     return cost
955 # Add more path constraints if needed
956
957 # repeat for more stages
958 ...
959 ...
960
961 Structure your output in a single python code block as follows for bimanual robot:
962 ```python
963
964 # Your explanation of how many stages are involved in the task and what each stage is about.
965 # ...
966
967 num_stages = ?
968
969 ### left-arm stage 1 sub-goal constraints (if any)
970 def left_stage1_subgoal_constraint1(end_effector, keypoints):
971     """Put your explanation here."""
972     ...
973     return cost
974
975 ### right-arm stage 1 sub-goal constraints (if any)
976 def right_stage1_subgoal_constraint1(end_effector, keypoints):
977     """Put your explanation here."""
978     ...
979     return cost
980 # Add more sub-goal constraints if needed
981
982 ### left stage 1 path constraints (if any)
983 def left_stage1_path_constraint1(end_effector, keypoints):
984     """Put your explanation here."""
985     ...
986     return cost
987 ### right stage 1 path constraints (if any)
988 def right_stage1_path_constraint1(end_effector, keypoints):

```

```

989     """Put your explanation here."""
990     ...
991     return cost
992 # Add more path constraints if needed
993
994 # repeat for more stages
995 ...
996 ...
997
998 ## Query
999 Query Task: "[INSTRUCTION]"
1000 Query Image: [IMAGE WITH KEYPOINTS]

```

1002 A.7 Implementation Details of Point Tracker

1003 We implement a simple point tracker following [121] based on DINOv2 (ViT-S14) [5] that leverages
1004 the fact that multiple RGB-D cameras are present and DINOv2 is efficient to run at a real-time
1005 frequency.

1006 At initialization, an array of 3D keypoint positions $\mathbf{k} \in \mathbb{R}^3$ are given. We first take the RGB-D
1007 captures from each present camera. For each RGB image, we obtain the pixel-wise DINOv2 features
1008 following the same procedure in Section A.5 and record their associated 3D world coordinates using
1009 calibrated cameras. For each 3D keypoint positions, we aggregate all the features from points that are
1010 within 2cm from all the cameras. The mean of the aggregated features is recorded as the reference
1011 feature for each keypoint, which is kept fixed throughout the task.

1012 After initialization, at each time step, we similarly obtain the pixel-wise features from DINOv2 from
1013 all cameras with their 3D world coordinates. To track the keypoints, we calculate cosine similarity
1014 between features across all pixels and the reference features. The top 100 matches are selected for
1015 each keypoint with a cutoff similarity of 0.6. We then reject outliers for the selected matches by
1016 calculating median deviation ($m = 2$). Additionally, as the tracked keypoints may oscillate in a
1017 small region, we apply a uniform filter with a window size of 10 in the end. The entire procedure
1018 runs at a fixed frequency of 20 Hz.

1019 Note that the implemented point tracker is a simplification from [121] for real-time tracking. We
1020 refer readers to [121] for more comprehensive discussion on using self-supervised vision models,
1021 such as DINOv2, for point tracking. Alternatively, more specialized point trackers can be used [141–
1022 148].

1023 A.8 Implementation Details of Sub-Goal Solver

1024 The sub-goal problems are implemented and solved using SciPy [128]. The decision variable is a
1025 single end-effector pose (position and Euler angles) in \mathbb{R}^6 for single-arm robots and two end-effector
1026 poses in \mathbb{R}^{12} for bimanual robot. The bounds for the position terms are the pre-defined workspace
1027 bounds, and the bounds for the rotation terms are that the half hemisphere where the end-effector
1028 faces down (due to the joint limits of the Franka arm, it is often likely to reach joint limit when an
1029 end-effector pose faces up). The decision variables are normalized to $[-1, 1]$ based on the bounds.
1030 For the first solving iteration, the initial guess is chosen to be the current end-effector pose. We
1031 use sampling-based global optimization Dual Annealing [129] in the first iteration to quickly search
1032 the full space, which is followed by a gradient-based local optimizer SLSQP [130] that refines the
1033 solution. The full procedure takes around 1 second for this iteration. In subsequent iterations, we
1034 use the solution from previous stage and only use local optimizer as it can quickly adjust to small
1035 changes. The optimization is cut off with a fixed time budget represented as number of objective
1036 function calls to keep the system running at a high frequency.

1037 We discuss the cost terms in the objective function below.

1038 **Constraint Violation:** We implement constraints as cost terms in the optimization problem, where
1039 the returned costs by the ReKep functions are multiplied with large weights.

1040 **Scene Collision Avoidance:** We use nvblox [149] with the PyTorch wrapper [58] to compute the
1041 ESDF of the scene in a separate node that runs at 20 Hz. The ESDF calculation aggregates the

1042 depth maps from all available cameras and excludes robot arms using cuRobo and any grasped rigid
1043 objects (tracked via a masked tracker model Cutie [136]). A collision voxel grid is then calculated
1044 using the ESDF and used by other modules in the system. In the sub-goal solver module, we first
1045 downsample the gripper points and the grasped object points to have a maximum of 30 points using
1046 farthest point sampling. Then we calculate the collision cost using the ESDF voxel grid with linear
1047 interpolation with a threshold of 15cm.

1048 **Reachability:** Since our decision variables are end-effector poses, which may not be always reach-
1049 able by the robot arms, especially in confined spaces, we need to add a cost term that encourages
1050 finding solutions with valid joint configurations. Therefore, we solve an IK problem in each iteration
1051 of the sub-goal solver using PyBullet [133] and use its residual as a proxy for reachability. We find
1052 that this takes around 40% of the time of the full objective function. Alternatively, one may solve
1053 the problem in joint space, which would ensure the solution is within the joint limits by enforcing
1054 the bounds. We find that this is inefficient with our Python-based implementation as we need to cal-
1055 culate forward kinematics for a magnitude of more times in the path solver, because the constraints
1056 are evaluated in the task space. To address this while ensuring efficiency, future works can consider
1057 using hardware-accelerated implementations to solve the problems in joint space [58].

1058 **Pose Regularization:** We also add a small cost that encourages the sub-goal to be close to the
1059 current end-effector pose.

1060 **Consistency:** Since the solver iteratively solves the problem at a high frequency and the noise from
1061 the perception pipeline may propagate to the solver, we find it useful to include a consistency cost
1062 that encourages the solution to be close to the previous solution.

1063 **(Dual-Arm only) Self-Collision Avoidance:** To avoid two arms collide with each other, we compute
1064 the pairwise distance between the two point sets, each including the gripper points and grasped
1065 object points.

1066 A.9 Implementation Details of Path Solver

1067 The path problems are implemented and solved using SciPy [128]. The number of decision variables
1068 is calculated based on the distance between the current end-effector pose and the target end-effector
1069 pose. Specifically, we define a fixed step size (20cm and 45 degree) and linearly approximate the
1070 desired number of “intermediate poses”, which are used as decision variables. As in the sub-goal
1071 problem, they are similarly represented using position and Euler angles with the same bounds. For
1072 the first solving iteration, the initial guess is chosen to be linear interpolation between the start and
1073 the target. We similarly use sampling-based global optimization followed by a gradient-based local
1074 optimizer in the first iteration and only use local optimizer in subsequent iterations. After we obtain
1075 the solution, represented as a number of intermediate poses, we fit a spline using the current pose,
1076 the intermediate poses, and the target pose, which are then densely sampled to be executed by the
1077 robot.

1078 In the objective function, we first unnormalize the decision variables and use piecewise linear in-
1079 terpolation to obtain a dense sequence of discrete poses to represent the path (referred to as “dense
1080 samples” below). A spline interpolation would be aligned with how we postprocess and execute the
1081 solution, but we find linear interpolation to be computationally more efficient. Below we discuss the
1082 individual cost terms in the objective function.

1083 **Constraint Violation:** Similar to that in the sub-goal problem, we check violation of the ReKep
1084 constraints for each dense sample along the path and penalize with large weights.

1085 **Scene Collision Avoidance:** The calculation is similar to the sub-goal problem, except that it is
1086 calculated for each dense sample. We ignore the collision calculation with a 5cm radius near the
1087 start and the target poses, as this tends to stabilize the solution when solved at a high frequency due
1088 to various real-world noises. We additionally add a table clearance cost that penalizes the path from
1089 penetrating the table (or the bottom of the workspace for the wheeled single-arm robot).

1090 **Path Length:** We approximate the path length using the dense samples by taking the sum of their
1091 differences. Shorter paths are encouraged.

1092 **Reachability:** We solve an IK problem for each intermediate pose inside the objective function as
1093 in the sub-goal problem. See the sub-goal solver section for more details.

1094 **Consistency:** As in the sub-goal problem, we encourage the solution to be close to the previous
1095 one. Specifically, we store the dense samples from the previous iteration. To calculate the solution
1096 consistency, we use the pairwise distance between the two sequences (treated as two sets) as an
1097 efficient proxy. Alternatively, Hausdorff distance can be used.

1098 **(Dual-Arm only) Self-Collision Avoidance:** We similarly compute self-collision avoidance for the
1099 dual-arm platform as in the sub-goal problem. We also use pairwise distance between the two
1100 sequences to efficiently calculate this cost.

1101 A.10 Comparisons with Prior Works on Visual Prompting for Manipulation

1102 There has been several concurrent works investigating the application of visual prompting of VLMs
1103 to robotic manipulation [99–101, 112, 125]. Below we summarized the differences to highlight the
1104 contributions of this work.

1105 **Task DoF:** In this work, we focus on challenging tasks that require 6 DoF (single arm) or 12 DoF
1106 (two arms) motions. However, this is not trivial for existing VLMs which operate on 2D images – as
1107 quoted from MOKA [100], “current VLMs are not capable of reliably predicting 6-DoF motions”
1108 and PIVOT [101], “generalizing to higher dimensional spaces such as rotation poses even additional
1109 challenges”. To tackle this, one key insight from ReKep is that VLMs only need to implicitly specify
1110 full 3D rotations by reasoning about keypoints in (x, y, z) *Cartesian coordinates*. After this, actual
1111 3D rotations are solved by high-precision and efficient numerical solvers, effectively sidestepping
1112 the challenge of explicitly predicting 3D rotations. As a result, the same formulation also naturally
1113 generalizes to controlling multiple arms.

1114 **High-Level Planning:** While many works also consider multi-stage tasks via an language-based
1115 task planners which are independent from their methods, our formulation takes inspiration from
1116 TAMP and organically integrates high-level task planning with low-level actions in a unified contin-
1117 uous mathematical program. As a result, the method can naturally consider *geometric dependencies*
1118 across stages and do so *at a real-time frequency*. When a failure occurs, it would backtrack to a
1119 previous stage in which its conditions can still be satisfied. For example, in the “pouring tea” task,
1120 the robot can only start tilting the teapot when the teapot spout is aligned with the cup opening.
1121 However, if the cup is being moved in the process, it should level the teapot and re-align with the
1122 cup. Or if the teapot is being taken from the gripper, it should instead re-grasp the teapot.

1123 **Low-Level Execution:** A common issue with using VLMs is that it is computationally expensive
1124 to run, hindering the high-frequency perception-action feedback loops often required for many ma-
1125 nipulation tasks. As a result, most of existing works either consider the open-loop settings where
1126 visual perception is only used in the beginning or only consider the tasks where slow execution is
1127 acceptable. Instead, our formulation natively supports a high-frequency perception-action loops by
1128 coupling VLMs with a point tracker, which effectively enables reactive behaviors via closed-loop
1129 execution despite leveraging very large foundation models.

1130 **Visual Prompting Methods:** We uniquely consider using visual prompting for code-generation,
1131 where code may contain arbitrary arithmetic operations on a set of keypoints via visual referring
1132 expressions. Although a single point is limiting to capture complex geometric structure, multiple
1133 points and their relations can even specify *vectors*, *surfaces*, *volumes*, and their *temporal dependen-*
1134 *cies*. While being conceptually simple, this offers a much higher degree of flexibility which can
1135 fully specify 6 DoF or even 12 DoF motions.

1136 A.11 Extended Discussions on Limitations

1137 Herein we present additional limitations of the existing system.

1138 **Prompting and Robustness:** Although we have demonstrated that existing VLMs possess rudimentary capabilities at specifying ReKep constraints, we have observed that when dealing with tasks that span many stages with several temporally dependent constraints (A.14), the VLMs lack enough robustness to obtain consistent success.

1142 **Task-Space Planning:** To enable efficient planning, in this work we only consider planning in the task space with the end-effector poses as decision variables. However, we have observed that in certain scenarios, it may be kinematically challenging for robots to achieve the optimized poses as the solver does not explicitly account for the kinematics of the robot. Planning in joint space can likely resolve the issue but we find it to be less computationally efficient for our tasks.

1147 **Articulated Object Manipulation:** In this work, we do not investigate tasks involving articulated objects, as we observed this requires advanced spatial reasoning capabilities that are beyond those of existing VLMs. However, the ReKep formulation may be extended to such tasks by representing different types of joints also by “relations of keypoints”. For example, ReKep constraints can be written to constrain certain keypoints to move only alongside a line (prismatic joints) or a curve (revolute joints). To extend to these scenarios, finetuning may be required as in [150–152].

1153 **Bimanual Coordination:** Although we demonstrate the application of ReKep to bimanual manipulation, we also identify several important limitations in this domain. Notably, the challenges can be roughly categorized into those pertaining to semantic reasoning of keypoint relations by the VLM and those pertaining to solving for bimanual motions by the optimization solver. For semantic reasoning, to achieve bimanual folding, the VLM needs to possess certain spatial knowledge about which steps should/can be performed together by both arms. For example, the bottom of a shirt often needs to be grasped by two hands, each at one corner, in order to fold it upwards to align with the collar. Another example in blanket folding is to recognize that the bottom-left corner should be aligned with the top-left corner and the bottom-right corner should be aligned with the top-right corner, as other matching may lead to self-collision. For optimization solver, as bimanual motion planning dramatically increases the search space of possible motions, which slows down the overall pipeline and more frequently produces less optimal behaviors.

1165 A.12 Simulation Experiments

1166 We additionally implement ReKep in OmniGibson [153] for the *Pour Tea* task. It is compared to a monolithic learning-based baseline based on the transformer architecture [154] adopted from RVT [155, 156]. The baseline is trained via imitation learning on 100 expert demonstrations, where demonstrations are from scripted policies using privileged simulation information. Success rates are averaged across 100 trials and reported below. Although the monolithic policy excels in training scenarios given its access to expert demonstrations, we observe that ReKep performs significantly stronger in unseen settings, and more importantly, without the need of expert demonstrations.

	Seen Poses	Unseen Poses	Unseen Objects
Monolithic Policy	0.93	0.31	0.14
ReKep (Zero-Shot)	0.75	0.68	0.72

1173 A.13 Comparisons of Visual Feature Extractors for Keypoint Proposal

1174 Herein we provide qualitative comparisons of different methods for keypoint proposal. We compare three pre-trained visual feature extractors, each of which represents a popular class of pre-training methods: DINOv2 [5] (self-supervised pre-training), CLIP [94] (vision-language contrastive pre-training), and ViT [157] (supervised pre-training). We also compare to a variant that does not use

1178 Segment Anything (SAM) [132] for its objectness prior. In Fig. 8, we show the extracted feature
1179 maps (projected to RGB space) and their clustered keypoints for three different scenes.

1180 We would like to note two important observations from the comparisons: 1) objectness prior given
1181 by SAM is critical to constrain the keypoint proposal on objects in the scene instead of on the
1182 background, and 2) while most visual foundation models can provide useful guidance, DINOv2
1183 produces sharper features that can better distinguish fine-grained regions of an object. The first
1184 observation can be clearly made by comparing the last column with other candidate methods. The
1185 second observation can be made by noting the following places: 1) the unique cyan color on the
1186 cup handle in the first scene, 2) the unique colors of the box panels in the second scene, and 3) the
1187 blue/green contrast of the top panel and the side panel in the third scene. Similarly, CLIP provides
1188 different features between different object parts, but the features are less sharp than those of DINOv2
1189 (color saturating from one part to the other). ViT, on the other hand, produces least distinguishable
1190 features between object parts, especially when texture is similar. In general, our observations are
1191 aligned with other works that also apply DINOv2 for its fine-grained object understanding [114,
1192 121, 123].



Figure 8: Comparisons of different methods for keypoint proposal.

1193 **A.14 Case Study of Long-Horizon Tasks**

1194 To stress test the presented system, we additionally perform a case study of a long-horizon bimanual
1195 task of preparing breakfast tray. A successful completion requires 10 stages: 1) grasp table cloth,
1196 2) place table cloth inside the tray, 3) grasp bread from the plate, 4) place bread on top of the table
1197 cloth, 5) grasp mug and teacup simultaneously with two arms, 6) align mug and teapot, 7) pour tea
1198 into the mug, 8) place the mug inside the tray, 9) grasp the two tray handles simultaneously with two
1199 arms, and 10) lift up the tray and hand it to the human operator.

1200 This task presents significant challenges for many components in the system. Specifically, we find
1201 that existing pipeline for keypoint proposal and constraint specification are incapable of generating
1202 the full set of correct keypoints and the full sequence of the correct ReKep constraints. Additionally,
1203 we observe that as a result of multiple present objects, our point tracker is incapable of consistently
1204 tracking all required keypoints in the scene. As a result, we resort to manually annotated keypoints
1205 and constraints, as well as keypoint detection at only the start of each stage instead of persistent
1206 keypoint tracking. With the above modifications, we find that the system can reasonably perform
1207 the task. The initial and final configurations are shown below. The solutions for each stage are
1208 shown on the next page. The video result is available at rekep-cori.github.io.



Figure 9: Initial configuration and successful final configuration of the preparing breakfast tray task.

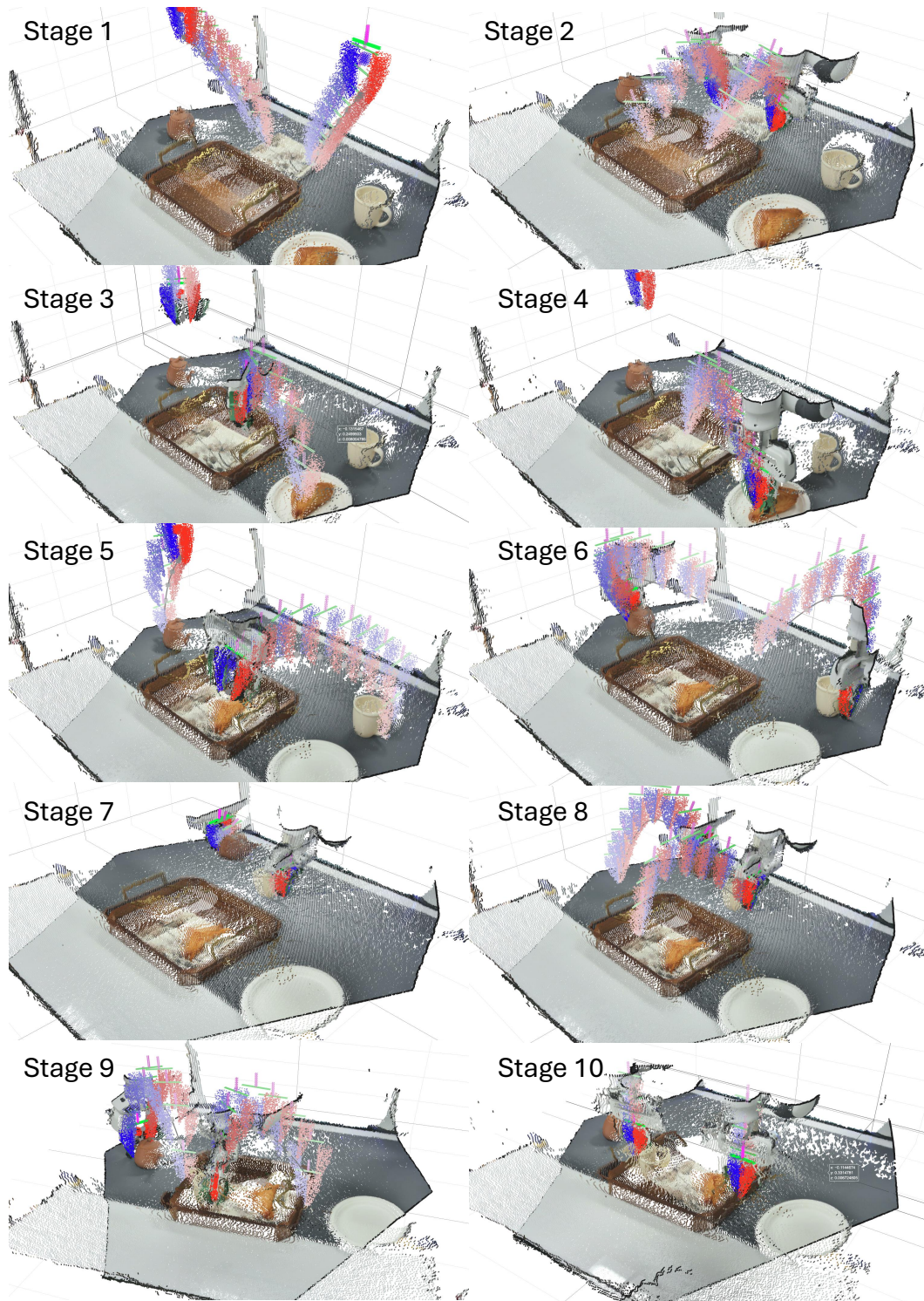


Figure 10: Solution visualization of the preparing breakfast tray task.