

Vision-based Manipulation from Single Human Video with Open-World Object Graphs

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** We present an object-centric approach to empower robots to learn
2 vision-based manipulation skills from human videos. We investigate the prob-
3 lem of imitating robot manipulation from a single human video in the *open-world*
4 setting, where a robot must learn to manipulate novel objects from one video
5 demonstration. We introduce ORION, an algorithm that tackles the problem by
6 extracting an object-centric manipulation plan from a single RGB-D video and de-
7 riving a policy that conditions on the extracted plan. ORION enables the robot to
8 learn from videos captured by daily mobile devices such as an iPad and generalize
9 the policies to deployment environments with varying visual backgrounds, cam-
10 era angles, spatial layouts, and novel object instances. We systematically evaluate
11 ORION on both short-horizon and long-horizon tasks, demonstrating the efficacy
12 of ORION in learning from a single human video in the open world.

13 **Keywords:** Robot Manipulation, Imitation From Human Videos

14 1 Introduction

15 A critical step toward building robot autonomy is developing sensorimotor skills for perceiving and
16 interacting with unstructured environments. Conventional methods for acquiring skills necessitate
17 manual engineering and/or costly data collection [1–5]. A promising alternative is teaching robots
18 through human videos of manipulation behaviors situated in everyday scenarios. These methods
19 have great potential to tap into the readily available source of Internet videos that encompass a wide
20 distribution of human activities, paving the ground for scaling up skill learning.

21 Prior work on learning from human videos has focused on pre-training representations and value
22 functions [6–10]. However, they do not explicitly capture object states and their interactions in 3D
23 space where robot motions are defined. Consequently, they require separate teleoperation data for
24 each set of objects in each location and even for each possible change in visual background, e.g.,
25 the scene background or lighting conditions [11]. In contrast, our goal is for a robot to imitate a task
26 robustly in the “open world”, i.e., under varying visual and spatial conditions from a single human
27 video, without prior knowledge of the object models or the behaviors shown. Since we consider
28 actionless videos that are equivalent to state-only demonstrations in the problem of “Imitation from
29 Observation”[12], we refer to our problem setting as *open-world imitation from observation*.

30 Developing a method in this setting is only possible due to the recent advances in vision foundation
31 models [13, 14]. These models, pre-trained on Internet-scale visual data, excel at understanding
32 open-vocabulary visual concepts and enable robots to recognize and localize objects in videos with-
33 out known object categories or access to physical states. This work marks the first step toward
34 achieving our vision of open-world imitation from observation, where a robot imitates how to in-
35 teract with objects given *a single video* while deployed in environments with different visual back-
36 grounds and unseen spatial configurations. In this work, we consider using RGB-D video demon-
37 strations where a person manipulates a small set of task-relevant objects with their single *hand*,

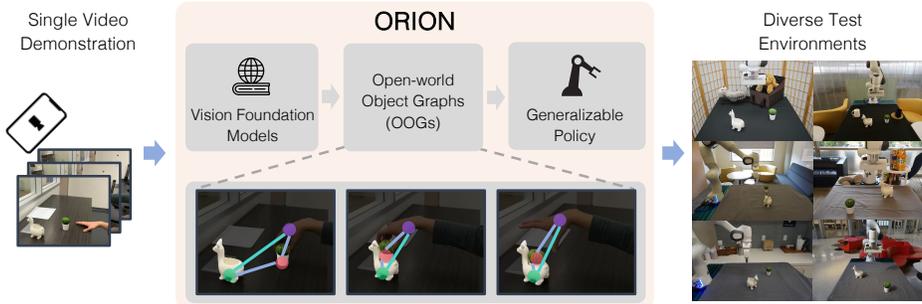


Figure 1: **Overview.** ORION tackles the problem of imitating manipulation from single human video demonstrations. ORION first extracts a sequence of Open-World Object Graphs (OOGs), where each OOG models a keyframe state with task-relevant objects and hand information. Then ORION leverages the OOG sequence to construct a manipulation policy that generalizes across varied initial conditions, specifically in four aspects: visual background, camera shifts, spatial layouts, and novel instances from the same object categories.

38 recorded with a stationary camera. These videos are actionless or state-only, as they do not come
 39 with any ground-truth action labels for the robot.

40 We introduce our method ORION, short for **Open-wORLD** video **Imitation**ON. Figure 1 visualizes a
 41 high-level overview of ORION. The core innovation lies in creating an object-centric spatiotemporal
 42 abstraction that effectively bridges the observational gap between human demonstration and robot
 43 execution. The design of ORION stems from our insight that manipulation tasks center around
 44 object interaction, and task completion depends on whether specific intermediate states, so-called
 45 *subgoals*, are reached. To capture the object-centric information in the video, we design a graph-
 46 based, object-centric representation, called Open-world Object Graphs (OOGs), to model the states
 47 of task-relevant objects and their relationships. An OOG has a two-level hierarchy. The high level
 48 consists of the object nodes and a hand node, where object nodes identify and localize the relevant
 49 objects by leveraging outputs from vision foundation models, while the hand node encodes the
 50 interaction information between the hand and objects, such as where to grasp. The low level consists
 51 of point nodes, which correspond to object keypoints, and the node features detail the motions of
 52 object keypoints in the 3D space.

53 ORION extracts a manipulation plan from the video as a sequence of OOGs and uses the plan to
 54 construct a generalizable policy. In experiments, ORION constructs a policy robust to conditions
 55 vastly different from the one in the video. Using only an iPhone or an iPad to record a human per-
 56 forming tasks in everyday environments (e.g., an office or a kitchen), ORION policies are deployed
 57 in workspaces with drastically different visual backgrounds, camera angles, and spatial arrange-
 58 ments, and even generalize to manipulating unseen object instances of the same categories.

59 In summary, our contribution is three-fold: 1) We pose the problem of learning vision-based robot
 60 manipulation from a single human video in the open-world setting; 2) We introduce Open-world
 61 Object Graphs (OOGs), a graph-based, object-centric representation for modeling the states and
 62 relations of task-relevant objects; and 3) We present ORION, an algorithm that uses a single video
 63 to construct a manipulation policy, which generalizes to conditions that differ in four key ways:
 64 visual backgrounds, camera perspectives, spatial configurations, and new object instances.

65 2 Problem Formulation

66 In this paper, we consider a vision-based, tabletop manipulation task, formulated as a finite-horizon
 67 Markov Decision Process (MDP) described by a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, H, R, \mu \rangle$, where \mathcal{S} is the state space
 68 of raw sensory data including RGB-D images and robot proprioception, \mathcal{A} is the action space of low-
 69 level robot commands, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$ is the transition dynamics, H is the maximal task horizon,
 70 R is the sparse reward function, and μ is the initial state distributions of a task. In this work, we
 71 consider the case where task reward functions are defined based on the contact relations between
 72 a small set of *task-relevant objects*. For example, a mug is placed on top of a coaster, or a spoon

73 is put inside a bowl. A reward function returns 1 if all object relations of a task are satisfied and 0
74 otherwise. The primary objective of solving a manipulation task is to find a visuomotor policy π that
75 maximizes the expected task success rate from a wide range of initial configurations, characterized
76 by μ , where the states vary across the following four dimensions: 1) changing visual backgrounds,
77 2) different camera angles, 3) different object instances from the same categories, and 4) varied
78 spatial layouts of the task-relevant objects.

79 We assume a robot does not have direct access to the ground-truth task reward or the physical
80 states of task-relevant objects. We consider a setting where a single *actionless* video [15, 16] V is
81 provided as a *state-only* demonstration. We assume V to be a video stream of a person manipulating
82 the task-relevant objects with their single *hand*, captured as a sequence of RGB-D images using
83 a stationary camera. V is an arbitrarily long video that involves a manipulation sequence where
84 the contact relations among task-relevant objects and the hand change (e.g., an object is grasped
85 or an object is placed on top of another). The assumption about V refers to tasks that involve
86 diverse manipulation behaviors such as pick-and-place, assembly, object insertion. To avoid the
87 inherent ambiguities of videos due to the distraction of irrelevant objects or ambiguities of what a
88 user wants to specify (whether the color of a task-relevant mug matters to the task or not), each
89 V is accompanied by a complete list of English descriptions of the task-relevant objects with their
90 complete feature descriptions such as their color that a user wants, uniquely defining the object
91 instances in V . Such a list is represented as a comma-separated list; an example is “[‘small red
92 block’, ‘boat body’]” for the task shown in Figure 2. In this scenario, however, the robot is not
93 pre-programmed to have access to ground-truth categories and locations of the task-relevant objects
94 in V . We refer to this challenging setting as “open-world” [17], as the robot must imitate from V
95 while not pre-programmed or trained to interact with the objects in V . To allow a robot to operate
96 in this “open-world” setting, we assume access to common sense knowledge through large models
97 pre-trained on internet-scale data, i.e., foundation models. For evaluation, we adopt the following
98 procedure. Given a single video V that accomplishes a task instance drawn from μ , the performance
99 of an approach is quantified by the average rewards received when evaluating new task instances
100 drawn from the same μ .

101 3 Method

102 We introduce **ORION (Open-world video ImitatiON)**, an algorithm that allows a robot to mimic
103 how to perform a manipulation task given a single human video, V . To effectively construct a
104 policy π from V , ORION employs a learning objective based on an object-centric prior. The goal
105 is to create a policy π that directs the robot to move objects along 3D trajectories that mimic the
106 directional and curvature patterns observed in V , relative to the objects’ initial and final positions.
107 This objective is based on the observation that objects are likely to achieve target configurations
108 by moving along trajectories similar to those in V . Key to ORION is generating a manipulation
109 plan from V , which serves as the spatiotemporal abstraction of the video that guides the robot to
110 perform a task. A plan is a sequence of object-centric keyframes that each specifies an initial or a
111 subgoal state captured in V . We first introduce our formulation of the object-centric representation
112 of a state, Open-world Object Graph (OOG), used in ORION, and then describe the algorithm that
113 constructs a robot policy given a human video.

114 3.1 Open-world Object Graph

115 At the core of our approach is a graph-based, object-centric representation, Open-world Object
116 Graphs (OOGs). OOGs use open-world vision models that model the visual scenes with task-
117 relevant objects and the hand such that they naturally exclude the distracting factors in visual data
118 and localize the task-relevant objects regardless of their spatial locations (see Section 3.2).

119 We denote an OOG as \mathcal{G} . At the high level, each object node corresponds to a task-relevant object
120 from the result of open-world vision models. Every object node comes with node features, consisting
121 of colored 3D point clouds derived from RGB-D observations. This node feature indicates both what

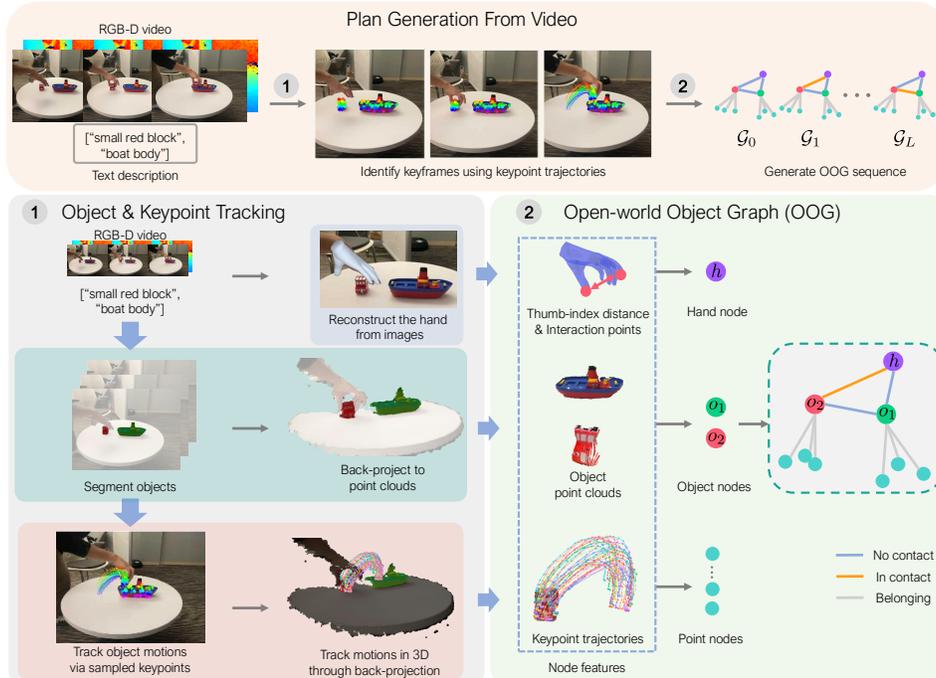


Figure 2: [Figure updated] **Overview of plan generation in ORION.** ORION generates a manipulation plan from a given video V in order for subsequent policies to synthesize actions. ORION first tracks the objects and keypoints across the video frames. Then keyframes are identified based on the velocity statistics of the keypoint trajectories. Then ORION generates an Open-world Object Graph (OOG) for every keyframe, resulting in a sequence of OOGs that serves as the spatiotemporal abstraction of the video. The figure is viewed best in color.

122 and where objects are and also represents their geometry information. Additionally, to inform the
 123 robot where to interact with objects (e.g. where to grasp), we introduce the specialized “hand node”,
 124 which stores the interaction cues such as contact points and the grip status (open or closed) that can
 125 be directly mapped to the robot end-effector during execution. At the low level, each point node
 126 corresponds to a keypoint that belongs to a task-relevant object. Every point node comes with the
 127 feature, namely the 3D motion trajectories. The feature explicitly models how an object should be
 128 moved during a manipulation task. In the rest of the paper, by motion features of a point node in G_l ,
 129 we mean 3D trajectory between keyframe l and $l + 1$.

130 In an OOG, all the object nodes and the hand node are fully connected, reflecting real-world spatial
 131 relationships. Each edge is augmented with a binary attribute that indicates if two objects or objects
 132 and the hand are in contact. This attribute allows our algorithm to check the set of satisfied contact
 133 relations, retrieving the matched OOG from the generated plan (see Section 3.2). The low-level point
 134 nodes are connected to their respective object node, indicating a belonging relationship. We denote
 135 node entities from human videos with a superscript V , and denote the ones from the robot rollout
 136 with a superscript Ro . Table 1 in the appendix also summarizes the variables needed in an OOG.

137 3.2 Manipulation Plan Generation From V

138 We describe the first part of ORION (see Figure 2), which automatically annotates the video and
 139 generates a manipulation plan from V . Here, a manipulation plan is a spatiotemporal abstraction
 140 of V that centers around the object states and their motions over time. Our core insight is that a
 141 task can be cost-effectively modeled with object locations at some keyframe states where the set of
 142 satisfied contact relations are changed, and abstract the rest of the states into 3D motions of objects.
 143 Concretely, a plan is represented as a sequence of OOGs, $\{G_l\}_{l=0}^L$ which corresponds to $L + 1$
 144 keyframes in V , with G_0 representing the initial state.

145 **Tracking task-relevant objects.** ORION first localizes task-relevant objects in the video V . Given
 146 V and the list of object descriptions mentioned in Section 2, ORION uses an open-world vision

147 model, Grounded-SAM [18], to annotate video frames with segmentation masks of the task-relevant
148 objects. In practice, due to the demanding computation of using open-world vision models, we
149 reduce the computation by exploiting object permanence to track the objects. Specifically, ORION
150 annotates the first video frame with Grounded-SAM, and then propagates the segmentation to the
151 rest of the video using a Video Object Segmentation Model, Cutie [19].

152 **Discovering keyframes.** After annotating the locations of task-relevant objects, we track their mo-
153 tions across the video to discover the keyframes based on the velocity statistics of object motions.
154 This design is based on the observation that changes in object contact relations due to manipu-
155 lation are often accompanied by sudden changes in object motions (e.g., transitioning from free
156 space motion to grasping an object). However, keeping full track of object point motion using tech-
157 niques like optical flow estimation requires heavy computation and the tracking quality is suscepti-
158 ble to noisy observations, largely due to occlusions during manipulation. We use a Track-Any-Point
159 (TAP) model, namely CoTracker [20], to track a subset of points in a long-term video with explicit
160 occlusion modeling, which has been successfully applied to track object motions in robot manipula-
161 tion [21, 22]. Specifically, we first sample keypoints within the object segmentation of the first frame
162 and track the trajectories across the video. The changes in velocity statistics are straightforward to
163 detect based on the TAP trajectories, where we discover the keyframes using a standard unsuper-
164 vised changepoint detection algorithm [23], a common technique that has been used in robotics
165 applications [24, 25].

166 **Generating OOGs from V .** Once ORION discovers the keyframes, it generates an OOG at
167 each keyframe to model the state of task-relevant objects and the human hand in V . The creation
168 of OOG nodes reuses the results from the annotation process: for object nodes, the point clouds
169 for node features are back-projected from the object segmentation using depth data; for the point
170 nodes, each node corresponds to the sampled keypoints, and their motion features, 3D trajectories,
171 are back-projected from the TAP trajectories using depth. Additionally, hand information is required
172 to specify the interaction points with task-relevant objects and the grip status to be mapped to the
173 robot gripper. We use a hand-reconstruction model, HaMeR [26], which gives a reconstructed hand
174 mesh that pinpoints the hand locations at each keyframe. The distances between the fingertips of the
175 mesh help determine the grip status, i.e., whether it is open or closed.

176 With all the node information, ORION establishes the edge connections between nodes in OOGs,
177 representing contact relations. Since all object and hand locations are computed in the camera frame
178 while the camera extrinsic of V is unknown, there is ambiguity when deciding the spatial rela-
179 tions between objects. We exploit the assumption of tabletop manipulation, where a table is always
180 present with its normal direction aligned with the z-axis of the world coordinate system. So ORION
181 estimates the transformation matrix of the table plane and transforms all the point cloud features in
182 OOGs to align with the xy plane of the world coordinate (Full details appear in Appendix C.2).
183 Then, the contact relations in each state can be determined based on the spatial relations and the
184 computed distances between point clouds. The relations allow ORION to match the test-time ob-
185 servations with a keyframe state from the plan and subsequently decide which object to manipulate
186 (see Section 3.3). In the end, ORION generates a complete OOG for each discovered keyframe.

187 3.3 Robot Policy To Synthesize Actions

188 Given a manipulation plan, ORION derives a manipulation policy that synthesizes actions based on
189 the aforementioned objective to achieve object motion similarities (detailed in Figure 3). The action
190 synthesis comprises three major steps: identify a keyframe from the plan that matches the current
191 observation, predict object motions, and use the predictions to optimize the robot actions. These
192 three steps are repeated until a task is completed or fails, detailed in Appendix E. The resulting
193 ORION policy is robust to visual variations due to the use of open-world vision models. It also
194 generalizes to different spatial locations due to our choice of representing object locations in object-
195 centric frames and the optimization process that is not constrained to specific positions.

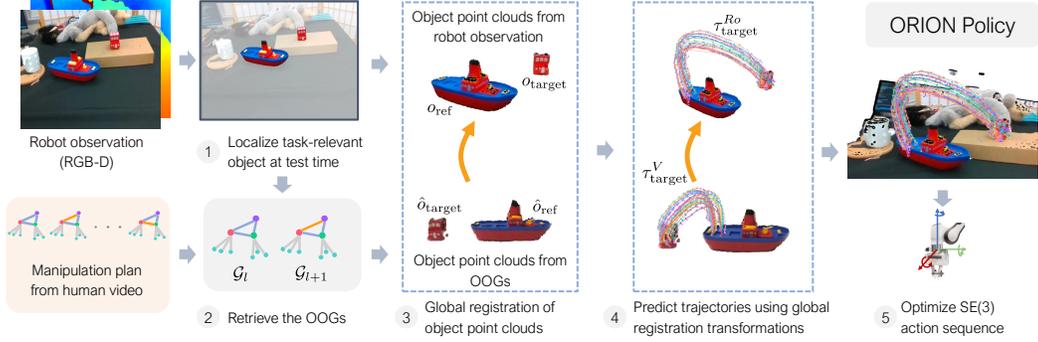


Figure 3: **Overview of the ORION Policy.** (1) ORION first localizes task-relevant objects at test time. (2) Next, ORION retrieves the matched OOGs from the generated manipulation plan. (3) ORION obtains the point clouds of the target object from the observation and the OOGs, namely o_{target} and \hat{o}_{target} , and those of the reference object, o_{ref} and \hat{o}_{ref} . Global registration is then performed to compute two transformations, one from \hat{o}_{target} to o_{target} , and the other from \hat{o}_{ref} to o_{ref} . (4) ORION then uses the computed transformations to warp τ_{target}^V , keypoint trajectories of the target object from the OOGs, into the workspace (details in the main text). The trajectory warping results in a predicted trajectory $\tau_{\text{target}}^{Ro}$. (5) ORION then uses $\tau_{\text{target}}^{Ro}$ to optimize the SE(3) action sequence of the robot end effector, which is subsequently used to command the robot.

196 **Retrieving OOGs from the plan.** ORION identifies the keyframe and retrieves OOGs to help
 197 decide what next actions to take. At test-time, ORION localizes objects in the new observations
 198 and estimates contact relations using the same vision pipeline as described in Section 3.2. Then
 199 ORION retrieves the OOG that has the same set of relations as the current state, allowing us to
 200 identify a pair $(\mathcal{G}_l, \mathcal{G}_{l+1})$, where \mathcal{G}_l is the retrieved graph and \mathcal{G}_{l+1} the graph of the next keyframe.
 201 This pair of graphs provides sufficient information to decide which object to manipulate next, termed
 202 the *target object*, and we denote its point cloud at keyframe l as \hat{o}_{target} , and its keypoint trajectories
 203 as τ_{target}^V . A target object is the one in motion due to manipulation between two keyframes, and it
 204 is determined by computing the average velocity per-object using motion features in \mathcal{G}_l . At the
 205 same time, another object, called the *reference object*, serves as a spatial reference for the target
 206 object’s movement when contact state relations change from \mathcal{G}_l to \mathcal{G}_{l+1} . We use the point cloud of
 207 the reference object at *next keyframe* $l + 1$, as object interactions might cause state changes of the
 208 reference object, and the information from the next keyframe gives us an accurate prediction of the
 209 trajectories. Once the target and reference objects are determined, we localize the corresponding
 210 objects in the new observations and denote their point clouds as o_{target} and o_{ref} , respectively.

211 **Predicting object motions.** Given the target and reference objects from keyframes l , and $l + 1$, we
 212 predict the motion of the target object in the current state by warping the keypoint trajectories esti-
 213 mated from V . To warp the trajectories, we first identify the initial and goal locations of keypoints in
 214 the new configuration by leveraging information given by the OOG pair. We use global registration
 215 of point clouds [27] to align \hat{o}_{target} with o_{target} and \hat{o}_{ref} with o_{ref} , giving us two transformations to
 216 compute the new starting and goal positions of target object keypoints conditioned on where the refer-
 217 ence object is. Then we normalize τ_{target}^V with its starting and goal locations, obtaining $\hat{\tau}_{\text{target}}^V$. $\hat{\tau}_{\text{target}}^V$
 218 only contains the directional and curvature patterns that are independent of the absolute location of
 219 the initial and the goal keypoints. Then we scale it back to the workspace coordinate frame using
 220 the new starting and goal locations, resulting in new keypoint trajectories of the target object $\tau_{\text{target}}^{Ro}$.

221 **Optimizing robot actions.** Once we obtain $\tau_{\text{target}}^{Ro}$, we optimize for a sequence of SE(3) transforma-
 222 tions that guide the robot end-effector to move. The SE(3) transformations are optimized to align
 223 the keypoint locations from previous frames to the next frames along the predicted trajectories:

$$224 \min_{T_0, T_1, \dots, T_{t_{l+1}-t_l}} \sum_{i=0}^{t_{l+1}-t_l} (\tau_{\text{target}}^{Ro}(i+1) - T_i \tau_{\text{target}}^{Ro}(i)) \quad (1)$$

224 where $\tau_{\text{target}}^{Ro}(i)$ ($0 \leq i \leq t_{l+1} - t_l$) represents the keypoint locations at timestep i along the trajec-
 225 tory. This optimization process naturally allows generalizations over spatial variations, as the action
 226 sequence always conditions on a new location instead of overfitting to fixed locations. To further
 227 specify where the gripper should interact with the object and whether it should be open or closed,



Figure 4: [Figure updated] This figure includes the following: task names, the initial and final frames of human videos, the list of word descriptions provided along with videos, snapshots of robot evaluation, and overall policy evaluation over all seven tasks, including the success rates and the quantification of failed trials, separated by failure mode. “Missed tracking” is the perception failure due to the vision foundation models, specifically the case of test-time object localization using Grounded-SAM.

228 we augment the resulting SE(3) sequence with the interaction information stored in the hand node
 229 h . To determine the initial pose of the end-effector in the sequence, ORION maps the two contact
 230 points using the computed transformation between \hat{o}_{target} and o_{target} . The mapped points correspond
 231 to the two finger tips of the robot gripper, and the robot’s gripper pose is determined by solving
 232 a simple inverse kinematics problem using the robot URDF file. We implement a combination of
 233 inverse kinematics (IK) and joint impedance control to achieve precise and compliant execution.

234 4 Experiments

235 In this section, we report on experiments to answer the following questions regarding the effec-
 236 tiveness of ORION and the important design choices. 1) Is ORION effective at constructing ma-
 237 nipulation policies given a single human video in the open-world setting? 2) To what extent does
 238 the object-centric abstraction improve the policy performance? 3) How critical is it to model the
 239 object motions with keypoints and the TAP formulation? 4) How consistent is the performance of
 240 ORION’s policy given videos taken in different conditions? 5) How effectively does ORION scale
 241 to long-horizon manipulation tasks?

242 4.1 Experiment Setup

243 **Task descriptions.** We design seven tasks to evaluate ORION poliiies: 1) Mug-on-coaster:
 244 placing a mug on the coaster; 2) Simple-boat-assembly: putting a small red
 245 block on a toy boat; 3) Chips-on-plate: placing a bag of chips on the plate;
 246 4) Succulents-in-llama-vase: inserting succulents into the llama vase; 5)
 247 Rearrange-mug-box: placing a mug on a coaster and placing a cream cheese box on a
 248 plate consecutively; 6) Complex-boat-assembly: placing both a small red block and a
 249 chimney-like part on top of a boat. 7) Prepare-breakfast: placing a mug on a coaster and
 250 putting a food box and can on the plate. The first four are “short-horizon” tasks that only require one
 251 contact relation between two objects, and the last three are “long-horizon” tasks that require more
 252 than one contact relation. Detailed success conditions of all tasks are described in Appendix E.
 253 Details about video recording, robot setup and evaluation can be found in Appendix B.

254 **Baselines.** To understand the model capacity and validate our design choices, we compare ORION
 255 with baselines. Since no prior work exists that matches the exact setting of our approach, we adopt
 256 the most important components from prior works and treat them as baselines to our model. Specifi-
 257 cally, we implement the following two baselines: 1) HAND-MOTION-IMITATION [9, 28] is a baseline

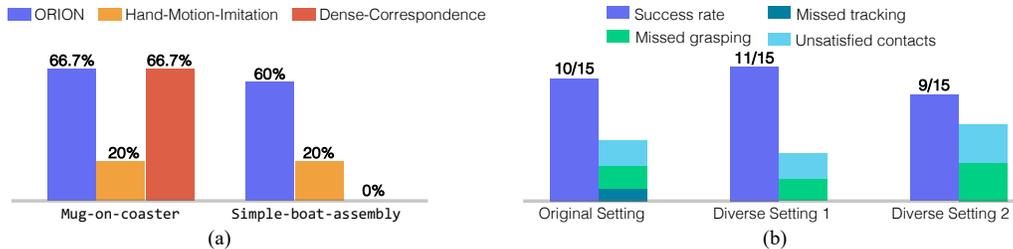


Figure 5: (a) Comparison experiments between ORION and the two baselines, namely HAND-MOTION-IMITATION and DENSE-CORRESPONDENCE. (b) Ablation study on using different videos of the task Mug-on-coaster. We show the number of successful trials out of 15 total trials on the bar plots for each setting. Figure 6 in Appendix F visualizes the different settings in this experiment.

258 that predicts robot actions by learning from the hand trajectories. The rest of the parts remain the
 259 same as ORION. We use this baseline to show whether it is critical to compute actions centering
 260 around objects. 2) DENSE-CORRESPONDENCE [15, 29] is a baseline that replace the TAP model
 261 in ORION with a dense correspondence model, optical flows. This baseline is used to evaluate
 262 whether our choice of TAP model is a better design. For this ablative study, we conduct experiments
 263 on Mug-on-coaster and Simple-boat-assembly to validate our model design, covering
 264 the distribution of common daily objects and assembly manipulation that requires precise control.

265 4.2 Experimental Results

266 Our evaluations are presented in Figures 4 and 5. We answer question (1) by showing the successful
 267 deployment of the ORION policies, while no other methods are designed to be able to operate in
 268 our setting. Furthermore, ORION yields an average of 66.7% success rates, which validates our
 269 model design in imitating from a single human video in the open-world setting.

270 We then answer question (2), showing the comparison results in Figure 5(a) against the baseline,
 271 HAND-MOTION-IMITATION, which yields low success rates in both tasks. Concretely, HAND-
 272 MOTION-IMITATION typically succeeds in trials where the initial spatial layouts are similar to the
 273 one in V . Its major failure mode is not being able to reach the target object configuration, e.g.,
 274 misplacing the mug on the table while not achieving contact with the coaster. These results imply
 275 that learning from human hand motion from V results in poor generalization abilities of policies,
 276 supporting the design choice of ORION which focuses on the object-centric information.

277 We further answer question (3) by comparing the performance between ORION and the op-
 278 tical flow baseline, DENSE-CORRESPONDENCE. The baseline performs drastically worse on
 279 Simple-boat-assembly than on Mug-on-coaster. Our further investigation shows that
 280 the baseline discovers keyframes in the middle of smooth transitions as opposed to changes in
 281 object contact relations, resulting in a manipulation plan that computes completely wrong actions.

282 To answer question (4), we conduct controlled experiments using the task Mug-on-coaster. We
 283 record two additional videos of the task in very different visual conditions and spatial layouts (see
 284 pictures in Appendix F) and construct a policy from each video. Then, we compare the two policies
 285 against the original one using the same set of evaluation conditions. The result in Figure 5(b) shows
 286 that there is no statistically significant difference in the performance, demonstrating that ORION is
 287 robust to videos taken under different visual conditions. Finally, we show that ORION is effective
 288 in scaling to long-horizon tasks. This conclusion is supported by the performance among the pairs
 289 of Mug-on-coaster versus Rearrange-mug-box, and Simple-boat-assembly versus
 290 Complex-boat-assembly. Both the short-horizon tasks are subgoals of their long-horizon
 291 counterparts, yet we do not see any performance drop between the two. Such result shows that
 292 ORION excels at scaling to long-horizon tasks without a significant drop in performance.

References

- 293
- 294 [1] M. Dalal, D. Pathak, and R. R. Salakhutdinov. Accelerating robotic reinforcement learning
295 via parameterized action primitives. *Advances in Neural Information Processing Systems*, 34:
296 21847–21859, 2021.
- 297 [2] A. Mandlekar, D. Xu, R. Martín-Martín, S. Savarese, and L. Fei-Fei. Learning to general-
298 ize across long-horizon tasks from human demonstrations. *arXiv preprint arXiv:2003.06085*,
299 2020.
- 300 [3] S. Nasiriany, H. Liu, and Y. Zhu. Augmenting reinforcement learning with behavior primitives
301 for diverse manipulation tasks. In *2022 International Conference on Robotics and Automation*
302 (*ICRA*), pages 7477–7484. IEEE, 2022.
- 303 [4] R. Zhang, S. Lee, M. Hwang, A. Hiranaka, C. Wang, W. Ai, J. J. R. Tan, S. Gupta, Y. Hao,
304 G. Levine, et al. Noir: Neural signal operated intelligent robots for everyday activities. *arXiv*
305 *preprint arXiv:2311.01454*, 2023.
- 306 [5] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. Viola: Imitation learning for vision-based manipulation
307 with object proposal priors. *arXiv preprint arXiv:2210.11339*, 2022.
- 308 [6] A. S. Chen, S. Nair, and C. Finn. Learning generalizable robotic reward functions from” in-
309 the-wild” human videos. *arXiv preprint arXiv:2103.16817*, 2021.
- 310 [7] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual represen-
311 tation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- 312 [8] Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang. Vip: Towards
313 universal visual reward and representation via value-implicit pre-training. *arXiv preprint*
314 *arXiv:2210.00030*, 2022.
- 315 [9] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mimicplay:
316 Long-horizon imitation learning by watching human play. *arXiv preprint arXiv:2302.12422*,
317 2023.
- 318 [10] H. Xiong, Q. Li, Y.-C. Chen, H. Bharadhwaj, S. Sinha, and A. Garg. Learning by watching:
319 Physical imitation of manipulation skills from human videos. In *2021 IEEE/RSJ International*
320 *Conference on Intelligent Robots and Systems (IROS)*, pages 7827–7834. IEEE, 2021.
- 321 [11] Y. Zhu, Z. Jiang, P. Stone, and Y. Zhu. Learning generalizable manipulation policies with
322 object-centric 3d representations. In *7th Annual Conference on Robot Learning*, 2023.
- 323 [12] F. Torabi. *Imitation Learning from Observation*. PhD thesis, University of Texas at Austin,
324 2021. PhD Thesis.
- 325 [13] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead,
326 A. C. Berg, W.-Y. Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- 327 [14] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haz-
328 iza, F. Massa, A. El-Nouby, et al. Dinov2: Learning robust visual features without supervision.
329 *arXiv preprint arXiv:2304.07193*, 2023.
- 330 [15] P.-C. Ko, J. Mao, Y. Du, S.-H. Sun, and J. B. Tenenbaum. Learning to act from actionless
331 videos through dense correspondences. *arXiv preprint arXiv:2310.08576*, 2023.
- 332 [16] C. Wen, X. Lin, J. So, K. Chen, Q. Dou, Y. Gao, and P. Abbeel. Any-point trajectory modeling
333 for policy learning. *arXiv preprint arXiv:2401.00025*, 2023.
- 334 [17] K. Joseph, S. Khan, F. S. Khan, and V. N. Balasubramanian. Towards open world object detec-
335 tion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
336 pages 5830–5840, 2021.

- 337 [18] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, et al. Grounding
338 dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint*
339 *arXiv:2303.05499*, 2023.
- 340 [19] H. K. Cheng, S. W. Oh, B. Price, J.-Y. Lee, and A. Schwing. Putting the object back into video
341 object segmentation. *arXiv preprint arXiv:2310.12982*, 2023.
- 342 [20] N. Karaev, I. Rocco, B. Graham, N. Neverova, A. Vedaldi, and C. Rupprecht. Cotracker: It is
343 better to track together. *arXiv preprint arXiv:2307.07635*, 2023.
- 344 [21] M. Vecerik, C. Doersch, Y. Yang, T. Davchev, Y. Aytar, G. Zhou, R. Hadsell, L. Agapito, and
345 J. Scholz. Robotap: Tracking arbitrary points for few-shot visual imitation. *arXiv preprint*
346 *arXiv:2308.15975*, 2023.
- 347 [22] B. Wen, W. Lian, K. Bekris, and S. Schaal. You only demonstrate once: Category-level ma-
348 nipulation from single visual demonstration. *arXiv preprint arXiv:2201.12716*, 2022.
- 349 [23] R. Killick, P. Fearnhead, and I. A. Eckley. Optimal detection of changepoints with a linear
350 computational cost. *Journal of the American Statistical Association*, 107(500):1590–1598,
351 2012.
- 352 [24] C. Wen, J. Lin, J. Qian, Y. Gao, and D. Jayaraman. Keyframe-focused visual imitation learning.
353 *arXiv preprint arXiv:2106.06452*, 2021.
- 354 [25] S. Niekum, S. Osentoski, C. G. Atkeson, and A. G. Barto. Online bayesian changepoint de-
355 tection for articulated motion models. In *2015 IEEE international conference on robotics and*
356 *automation (ICRA)*, pages 1468–1475. IEEE, 2015.
- 357 [26] G. Pavlakos, D. Shan, I. Radosavovic, A. Kanazawa, D. Fouhey, and J. Malik. Reconstructing
358 hands in 3d with transformers. *arXiv preprint arXiv:2312.05251*, 2023.
- 359 [27] S. Choi, Q.-Y. Zhou, and V. Koltun. Robust reconstruction of indoor scenes. In *Proceedings*
360 *of the IEEE conference on computer vision and pattern recognition*, pages 5556–5565, 2015.
- 361 [28] H. Bharadhwaj, A. Gupta, S. Tulsiani, and V. Kumar. Zero-shot robot manipulation from
362 passive human videos. *arXiv preprint arXiv:2302.02011*, 2023.
- 363 [29] N. Heppert, M. Argus, T. Welschehold, T. Brox, and A. Valada. Ditto: Demonstration imitation
364 by trajectory transformation. *arXiv preprint arXiv:2403.15203*, 2024.
- 365 [30] S. Bahl, A. Gupta, and D. Pathak. Human-to-robot imitation in the wild. *arXiv preprint*
366 *arXiv:2207.09450*, 2022.
- 367 [31] S. Kumar, J. Zamora, N. Hansen, R. Jangir, and X. Wang. Graph inverse reinforcement learning
368 from diverse videos. In *Conference on Robot Learning*, pages 55–66. PMLR, 2023.
- 369 [32] Y. Liu, A. Gupta, P. Abbeel, and S. Levine. Imitation from observation: Learning to imitate
370 behaviors from raw video via context translation. In *2018 IEEE International Conference on*
371 *Robotics and Automation (ICRA)*, pages 1118–1125. IEEE, 2018.
- 372 [33] P. Sharma, D. Pathak, and A. Gupta. Third-person visual imitation learning via decoupled
373 hierarchical controller. *Advances in Neural Information Processing Systems*, 32, 2019.
- 374 [34] L. Smith, N. Dhawan, M. Zhang, P. Abbeel, and S. Levine. Avid: Learning multi-stage tasks
375 via pixel-level translation of human videos. *arXiv preprint arXiv:1912.04443*, 2019.
- 376 [35] M. Xu, Z. Xu, C. Chi, M. Veloso, and S. Song. Xskill: Cross embodiment skill discovery. In
377 *Conference on Robot Learning*, pages 3536–3555. PMLR, 2023.
- 378 [36] R. Lee, J. Abou-Chakra, F. Zhang, and P. Corke. Learning fabric manipulation in the real world
379 with human videos. *arXiv preprint arXiv:2211.02832*, 2022.

- 380 [37] K. Shaw, S. Bahl, A. Sivakumar, A. Kannan, and D. Pathak. Learning dexterity from human
381 hand motion in internet videos. *The International Journal of Robotics Research*, 43(4):513–
382 532, 2024.
- 383 [38] A. Bahety, P. Mandikal, B. Abbatematteo, and R. Martín-Martín. Screwmimic: Bimanual
384 imitation from human videos with screw space projection. *arXiv preprint arXiv:2405.03666*,
385 2024.
- 386 [39] H. Bharadhwaj, A. Gupta, V. Kumar, and S. Tulsiani. Towards generalizable zero-shot manip-
387 ulation via translating human interaction plans. *arXiv preprint arXiv:2312.00775*, 2023.
- 388 [40] B. S. Pavse, F. Torabi, J. Hanna, G. Warnell, and P. Stone. Ridm: Reinforced inverse dynamics
389 modeling for learning from a single observed demonstration. *IEEE Robotics and Automation*
390 *Letters*, 5(4):6262–6269, 2020.
- 391 [41] H. Karnan, F. Torabi, G. Warnell, and P. Stone. Adversarial imitation learning from video
392 using a state observer. In *2022 International Conference on Robotics and Automation (ICRA)*,
393 pages 2452–2458. IEEE, 2022.
- 394 [42] F. Torabi, G. Warnell, and P. Stone. Imitation learning from video by leveraging proprioception.
395 In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages
396 3585–3591, 2019.
- 397 [43] F. Torabi, G. Warnell, and P. Stone. Generative adversarial imitation from observation. In
398 *Imitation, Intent, and Interaction (I3) Workshop at ICML 2019*, June 2019.
- 399 [44] F. Torabi, G. Warnell, and P. Stone. Behavioral cloning from observation. In *Proceedings of*
400 *the 27th International Joint Conference on Artificial Intelligence*, pages 4950–4957, 2018.
- 401 [45] Y. Duan, M. Andrychowicz, B. Stadie, O. Jonathan Ho, J. Schneider, I. Sutskever, P. Abbeel,
402 and W. Zaremba. One-shot imitation learning. *Advances in neural information processing*
403 *systems*, 30, 2017.
- 404 [46] N. Di Palo and E. Johns. Learning multi-stage tasks with one demonstration via self-replay. In
405 *Conference on Robot Learning*, pages 1180–1189. PMLR, 2022.
- 406 [47] S. Haldar, V. Mathur, D. Yarats, and L. Pinto. Watch and match: Supercharging imitation with
407 regularized optimal transport. In *Conference on Robot Learning*, pages 32–43. PMLR, 2023.
- 408 [48] S. Haldar, J. Pari, A. Rai, and L. Pinto. Teach a robot to fish: Versatile imitation from one
409 minute of demonstrations. *arXiv preprint arXiv:2303.01497*, 2023.
- 410 [49] E. Johns. Coarse-to-fine imitation learning: Robot manipulation from a single demonstration.
411 In *2021 IEEE international conference on robotics and automation (ICRA)*, pages 4613–4619.
412 IEEE, 2021.
- 413 [50] E. Valassakis, G. Papagiannis, N. Di Palo, and E. Johns. Demonstrate once, imitate imme-
414 diately (dome): Learning visual servoing for one-shot imitation learning. In *2022 IEEE/RSJ*
415 *International Conference on Intelligent Robots and Systems (IROS)*, pages 8614–8621. IEEE,
416 2022.
- 417 [51] A. Jonnavittula, S. Parekh, and D. P. Losey. View: Visual imitation learning with waypoints.
418 *arXiv preprint arXiv:2404.17906*, 2024.
- 419 [52] N. Di Palo and E. Johns. Dinobot: Robot manipulation via retrieval and alignment with vision
420 foundation models. *arXiv preprint arXiv:2402.13181*, 2024.
- 421 [53] D. Guo. Learning multi-step manipulation tasks from a single human demonstration. *arXiv*
422 *preprint arXiv:2312.15346*, 2023.

- 423 [54] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield. Deep object pose es-
424 timation for semantic robotic grasping of household objects. *arXiv preprint arXiv:1809.10790*,
425 2018.
- 426 [55] S. Tyree, J. Tremblay, T. To, J. Cheng, T. Mosier, J. Smith, and S. Birchfield. 6-dof pose esti-
427 mation of household objects for robotic manipulation: An accessible dataset and benchmark.
428 *arXiv preprint arXiv:2203.05701*, 2022.
- 429 [56] T. Migimatsu and J. Bohg. Object-centric task and motion planning in dynamic environments.
430 *IEEE Robotics and Automation Letters*, 5(2):844–851, 2020.
- 431 [57] D. Wang, C. Devin, Q.-Z. Cai, F. Yu, and T. Darrell. Deep object-centric policies for au-
432 tonomous driving. In *2019 International Conference on Robotics and Automation (ICRA)*,
433 pages 8853–8859. IEEE, 2019.
- 434 [58] C. Devin, P. Abbeel, T. Darrell, and S. Levine. Deep object-centric representations for gener-
435 alizable robot learning. In *2018 IEEE International Conference on Robotics and Automation*
436 *(ICRA)*, pages 7111–7118. IEEE, 2018.
- 437 [59] J. Shi, J. Qian, Y. J. Ma, and D. Jayaraman. Plug-and-play object-centric representations from
438 “what” and “where” foundation models.
- 439 [60] N. Di Palo and E. Johns. Keypoint action tokens enable in-context imitation learning in
440 robotics. *arXiv preprint arXiv:2403.19578*, 2024.
- 441 [61] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, B. Zitkovich,
442 F. Xia, C. Finn, et al. Open-world object manipulation using pre-trained vision-language mod-
443 els. *arXiv preprint arXiv:2303.00905*, 2023.
- 444 [62] K. Mo, P. Guerrero, L. Yi, H. Su, P. Wonka, N. Mitra, and L. J. Guibas. Structurenet: Hierar-
445 chical graph networks for 3d shape generation. *arXiv preprint arXiv:1908.00575*, 2019.
- 446 [63] Y. Huang, A. Conkey, and T. Hermans. Planning for multi-object manipulation with graph
447 neural network relational classifiers. In *2023 IEEE International Conference on Robotics and*
448 *Automation (ICRA)*, pages 1822–1829. IEEE, 2023.
- 449 [64] A. H. Qureshi, A. Mousavian, C. Paxton, M. C. Yip, and D. Fox. Nerp: Neural rearrangement
450 planning for unknown objects. *arXiv preprint arXiv:2106.01352*, 2021.
- 451 [65] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu. Hierarchical planning for long-horizon manip-
452 ulation with geometric and symbolic scene graphs. In *2021 IEEE International Conference on*
453 *Robotics and Automation (ICRA)*, pages 6541–6548. IEEE, 2021.
- 454 [66] M. Sieb, Z. Xian, A. Huang, O. Kroemer, and K. Fragkiadaki. Graph-structured visual imita-
455 tion. In *Conference on Robot Learning*, pages 979–989. PMLR, 2020.
- 456 [67] C. Truong, L. Oudre, and N. Vayatis. Selective review of offline change point detection meth-
457 ods. *Signal Processing*, 167:107299, 2020.
- 458 [68] Q.-Y. Zhou, J. Park, and V. Koltun. Open3d: A modern library for 3d data processing. *arXiv*
459 *preprint arXiv:1801.09847*, 2018.
- 460 [69] R. B. Rusu, N. Blodow, and M. Beetz. Fast point feature histograms (fpfh) for 3d registration.
461 In *2009 IEEE international conference on robotics and automation*, pages 3212–3217. IEEE,
462 2009.
- 463 [70] Z. Teed and J. Deng. Tangent space backpropagation for 3d transformation groups. In *Proceed-*
464 *ings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10338–
465 10347, 2021.

466 A Related Work

467 **Learning Manipulation From Human Videos.** Human videos offer a rich repertoire of object in-
468 teraction behaviors, making them an invaluable data source for manipulation. A large body of work
469 has explored how to leverage human video data for learning robot manipulation [9, 10, 30–35], either
470 through pre-training a single latent representation [7, 9, 35], learning representations of perception
471 or action priors [36, 37], learning an implicit reward function [6, 8], learning 6D representations of
472 actions [38], or learning generative models that in-paint human morphologies [15, 28, 30, 39]. How-
473 ever, they either require additional robot data from the target tasks or paired data between humans
474 and robots. Our approach takes a novel direction by tackling how a robot can imitate or learn from
475 a single human video only: the robot does not rely on pre-existing data, models, or ground-truth an-
476 notations *in scenes* where video recording and robot evaluation take place. We refer to such a setting
477 as *open-world imitation from observation*, where the robot is not programmed or trained to inter-
478 act with the objects in the video *a priori* and the video data does not come with any robot actions.
479 Our setting is closely related to the problem of “Imitation Learning from Observation” [12], where
480 state-only demonstrations are used to construct policies for physical interaction. However, this line
481 of prior work assumes simulators of demonstrated tasks exist and physical states of the agents or
482 objects are known [40–44]. In contrast, our setting does not assume the digital replica of real-world
483 tasks, and all the object information is only perceived through RGB-D videos.

484 **Learning Manipulation From a Single Demonstration.** Studies have delved into learning manip-
485 ulation policies from one demonstration. A notable one is one-shot imitation learning within meta-
486 learning framework proposed by Duan et al. [45]. While prior works on one-shot imitation learning
487 have shown a robot performing new tasks from one demonstration, they require extensive in-domain
488 data and a well-curated set of meta-training tasks beforehand, leading to significant data collection
489 costs and restricted policy generalization at test time due to the tailored nature of the training.

490 An alternative approach involves using a single demonstration for initial guidance, refining the pol-
491 icy through real-world self-play [22, 46–51]. However, this approach mainly applies to reset-free
492 tasks and struggles with scaling to multi-stage tasks where resetting to the task initial conditions
493 does not come free. Recently, foundation models are used to enable learning manipulation from
494 a single demonstration, but existing works require ground-truth access to the robot action through
495 kinesthetic teaching [52].

496 Our work aligns with these studies in using a single demonstration for learning manipulation, but
497 stands out by not needing prior data or self-play. Recent or concurrent works have also explored
498 using a single video demonstration only [29, 53], but they either assume known object instances
499 or lack in formulating systematic generalization in an open-world setting described in Section 2.
500 With just one single human video, our method constructs a policy that successfully completes the
501 task, while adapting to a wide range of visual and spatial differences from the task instance of video
502 demonstration.

503 **Object-Centric Representation for Learning Robot Manipulation.** The concept of object-
504 centric representation has long been recognized for its potential to enhance robotic perception and
505 manipulation by focusing on the objects within a scene. Prior works have shown effectiveness of
506 such representation in downstream manipulation tasks by factorizing visual scenes into disentangled
507 object concepts [54–58], but these works are typically confined to known object categories or
508 instances. Recent developments in foundation models allow robots to access the open-world
509 object concepts through pre-trained vision models [13, 14], enabling a wide range of abilities
510 such as imitation of long-horizon tabletop manipulation [5, 59], in-context learning of tabletop
511 manipulation [60], or mobile manipulation in the wild [61]. Building upon these advances, our work
512 focuses on leveraging open-world, object-centric concepts in imitating manipulation behaviors
513 from actionless human videos. We propose a graph-based representation called Open-world Object
514 Graph (OOG), which allows a robot to imitate from a human video by leveraging the object-centric
515 concepts. This proposed representation shares a similar vein with prior works that factorize scene or
516 task-relevant visual concepts into scene graphs [31, 62–66]. However, our representation is tailored

517 to integrate open-world object concepts and enable generalization across different embodiments,
 518 specifically a human and a robot.

519 B Additional Details of Experimental Setup

520 **Experimental setup.** We design experiments to fully test the efficacy of our method by providing
 521 the robot with videos captured in everyday scenarios, which naturally encompass visual back-
 522 grounds and camera setups that are different from the one for the robot. Specifically, we record
 523 an RGB-D video of a person performing each of the seven tasks in everyday scenarios, such as an
 524 office or a kitchen. We use an iPad for recording, which comes with a TrueDepth Camera, and we
 525 fix it on a camera stand. The videos can be found in the supplementary materials. During test time,
 526 the robot receives visual data through a single RGB-D camera, Intel Realsense435, and performs
 527 manipulation in its workstation to evaluate policies. We use the 7DoF Franka Emika Panda robot
 528 for all the experiments.

529 **Evaluation protocol.** As we describe in the experimental setup, the videos naturally include
 530 various visual backgrounds and camera perspectives that are significantly different from the robot
 531 workspace. Therefore, we only intentionally vary two dimensions before evaluating each trial of
 532 robot execution, namely the spatial layouts and the new object instances. Furthermore, the new
 533 object generalizations are included in the tasks *Mug-on-coaster* and *Chips-on-plate* as
 534 mugs and chip bags have many similar instances. As for the other three tasks, there are no novel
 535 objects involved, but we extensively vary the spatial layouts of task-relevant objects for evaluation.
 536 The policy performance of a task is the averaged success rates over 15 real-world trials. Aside from
 537 the success rates, we also group the failed executions into three types: *Missed tracking* of objects
 538 due to failure of the vision models, *Missed grasping* of objects during execution, and *Unsatisfied*
 539 *contacts* where the target object configurations are not achieved for reasons other than the previous
 540 two failure types.

541 C Additional Technical Details

542 C.1 Data Structure of an OOG.

543 For easy reproducibility of the proposed method, we present a table that explains the data structure
 544 of an OOG.

Node/Edge	Type	Attributes
$\mathcal{G}.vo_i$	Object Node	3D point cloud of an object.
$\mathcal{G}.vh$	Hand Node	Hand mesh and locations of the thumb and index finger.
$\mathcal{G}.vp_{ij}$	Point Node	A trajectory of a TAP keypoint between two keyframes, recorded in xyz positions.
$\mathcal{G}.eo_{ik}$	Object-Object Edge	A binary value of contact or not.
$\mathcal{G}.eh_i$	Object-Hand Edge	A binary value of contact or not.
$\mathcal{G}.ep_{ij}$	Object-Point Edge	The presence of an edge represents the belonging relation, and no specific feature is attached.

Table 1: Data Structure of an OOG. For a given OOG $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, it has $\mathcal{V} = \{\mathcal{G}.vo_i\} \cup \{\mathcal{G}.vh\} \cup \{\mathcal{G}.vp_{ij}\}$, and $\mathcal{E} = \{\mathcal{G}.eo_{ik}\} \cup \{\mathcal{G}.eh_i\} \cup \{\mathcal{G}.ep_{ij}\}$.

545 C.2 Implementation Details

546 **Changepoint detections.** We use changepoint detection to identify changes in velocity statistics of
 547 TAP keypoints. Specifically, we use a kernel-based changepoint detection method and choose radial

548 basis function [23]. The implementation of this function is directly based on an existing library
549 Ruptures [67].

550 **Plane estimation.** In Section 3.2, we mentioned using the prior knowledge of tabletop manipula-
551 tion scenarios and transforming the point clouds by estimating the table plane. Here, we explain
552 how the plane estimation is computed. Concretely, we rely on the plane estimation function from
553 Open3D [68], which gives an equation in the form of $ax + by + cz = d$. From this estimated
554 plane equation, we can infer a normal vector of the estimated table plane, (a, b, c) , in the camera
555 coordinate frame. Then, we align this plane with xy plane in the world coordinate frame, where we
556 compute a transformation matrix that displaces the normal vector (a, b, c) to the normalized vector
557 $(0, 0, 1)$ along the z-axis of the world coordinate frame. This transformation matrix is used to trans-
558 form point clouds in every frame so that the plane of the table always aligns with the xy plane of the
559 world coordinate.

560 **Object localization at test time.** When we localize objects at test time, there could be some false
561 positive segmentation of distracting objects. Such vision failures will prevent the robot policy from
562 successfully executing actions. To exclude such false positive object segmentation, we use Segmen-
563 tation Correspondence Model (SCM) from GROOT [11], where SCM filters out the false positive
564 segmentation of the objects by computing the affinity scores between masks using DINOv2 features.

565 **Global registration.** In this paper, we use global registration to compute the transformation between
566 observed object point clouds from videos and those from rollout settings. We implement this part
567 using a RANSAC-based registration function from Open3D [68]. Specifically, given two object
568 point clouds, we first compute their features using Fast-Point Feature Histograms (FPFH) [69], and
569 then perform a global RANSAC registration on the FPFH features of the point clouds [27].

570 **Implementation of SE(3) optimization.** We parameterize each homogeneous matrix T_i into a
571 translation variable and a rotation variable and randomly initialize each variable using the normal
572 distribution. We choose quaternions as the representation for rotation variables, and we normalize
573 the randomly initialized vectors for rotation so that they remain unit quaternions. With such param-
574 eterization, we optimize the SE(3) end-effector trajectories $T_0, T_1, \dots, T_{t_{i+1}-t_i}$ over the Objective
575 (1). However, jointly optimizing both translation and rotation from scratch typically results in trivial
576 solutions, where the rotation variables do not change much from the initialization due to the vanish-
577 ing gradients. To avoid trivial solutions, we implement a two-stage process. In the first stage, we
578 only optimize the rotation variables with 200 gradient steps. Then, the optimization proceeds to the
579 second stage, where we optimize both the rotation and translation variables for another 200 gradient
580 steps. In this case, we prevent the optimization process from getting stuck in trivial solutions for
581 rotation variables. We implement the optimization process using Lietorch [70].

582 D System Setup

583 **Details of camera observations.** As mentioned in Section 4, we use an iPad with a TrueDepth
584 camera for collecting human video demonstrations. We use an iOS app, Record3D, that allows us
585 to access the depth images from the TrueDepth camera. We record RGB and depth image frames
586 in sizes 1920×1080 and 640×480 , respectively. To align the RGB images with the depth data,
587 we resize the RGB frames to the size 640×480 . The app also automatically records the camera
588 intrinsics of the iPhone camera so that the back-projection of point clouds is made possible.

589 To stream images at test time, we use an Intel Realsense D435i. In our robot experiments, we use
590 RGB and depth images in the size 640×480 or 1280×720 in varied scenarios, all covered in our
591 evaluations. Evaluating on different image sizes showcases that our method is not tailored to specific
592 camera configurations, supporting the wide applicability of constructed policy.

593 **Implementation of real robot control.** In our evaluation, we reset the robot to a default joint
594 position before object interaction every time. Then we use a reaching primitive for the robot to reach
595 the interaction points. Resetting to the default joint position enables an unoccluded observation of
596 task-relevant objects at the start of each decision-making step. Note that the execution of object

597 interaction does not necessarily require resetting. To command the robot to interact with objects,
598 we convert the optimized SE(3) action sequence to a sequence of joint configurations using inverse
599 kinematics and control the robot using joint impedance control. We use the implementation of
600 Deoxys [5] for the joint impedance controller that operates at 500 Hz. To avoid abrupt motion and
601 make sure the actions are smooth, we further interpolate the joint sequence from the result of inverse
602 kinematics. Specifically, we choose the interpolation so that the maximal displacement for each joint
603 does not exceed 0.5 radian between two adjacent waypoints.

604 **E Success conditions of tasks**

605 We describe the success conditions for each of the tasks in detail:

- 606 • `Mug-on-coaster`: A mug is placed upright on the coaster.
- 607 • `Simple-boat-assembly`: A red block is placed in the slot closest to the back of the
608 boat. The block needs to be upright in the slot.
- 609 • `Chips-on-plate`: A bag of chips is placed on the plate, and the bag does not touch the
610 table.
- 611 • `Succulents-in-llama-vase`: A pot of succulents is inserted into a white vase in the
612 shape of a llama.
- 613 • `Rearrange-mug-box`: The mug is placed upright on the coaster, and the cream cheese
614 box is placed on the plate.
- 615 • `Complex-boat-assembly`: The chimney-like part is placed in the slot closest to the
616 front of the boat. The red block is placed in the slot closest to the back of the boat. Both
617 blocks need to be upright in the slots.
- 618 • `Prepare-breakfast`: The mug is placed on top of a coaster, the cream cheese box is
619 placed in the large area of the plate, and the food can is placed on the small area as shown
620 in the video demonstration.

621 In practice, we record the success and failure of a rollout as follows: If the program in ORION
622 policy returns true when matching the observed state with the final OOG from a plan, we mark a
623 trial as success as long as we observe that the object state indeed satisfies the success condition of
624 a task as described above. Otherwise, if the robot generates dangerous actions (bumping into the
625 table) or does not achieve the desired subgoal after executing the computed trajectory, we consider
626 the rollout as a failure and we manually record the failure.

627 **F Additional Details on Experiments**

628 **Diverse video recordings used in the ablation study.** Figure 6 shows the three videos taken in very
629 different scenarios: kitchen, office, and outdoor. The video taken in kitchen scenario is used in the
630 major quantitative evaluation, termed “Original setting”. The other two settings are termed “Diverse
631 setting 1” and “Diverse setting 2.” We conduct an ablation study where we compare policies imitated
632 from these three videos, which inherently involve varied visual scenes, camera perspectives. The
633 result of the ablation study is shown in Figure 5.

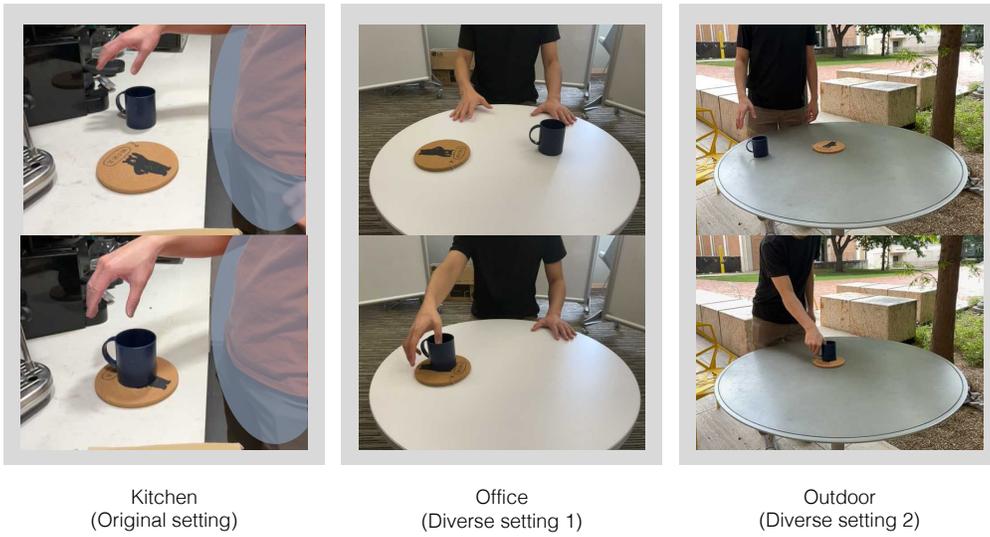


Figure 6: This figure visualizes the initial and final frames of the three videos of the same task *Mug-on-coaster*.