

Transformers are adaptable task planners

Anonymous Author(s)

Affiliation

Address

email

1 **Code:** https://anonymous.4open.science/r/temporal_task_planner-Paper148/

2 **A Hardware Experiments**

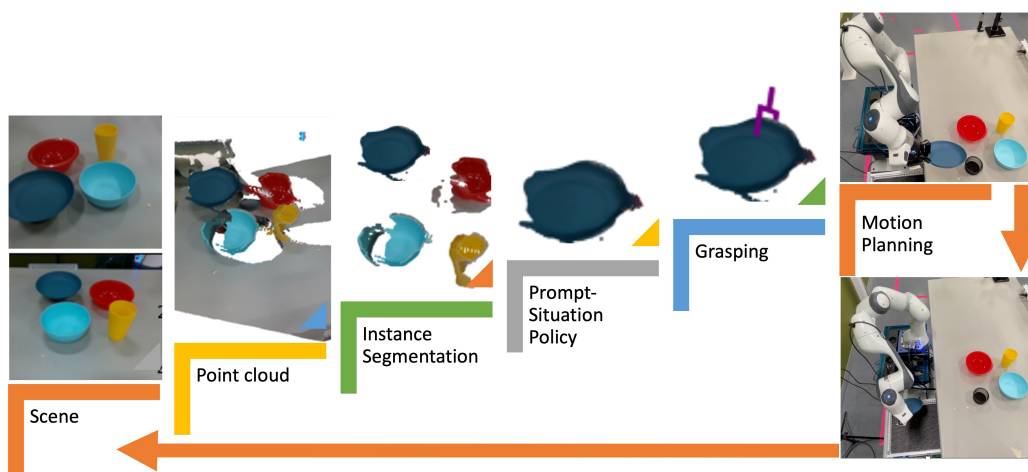


Figure 1: Pipeline for Real Hardware Experiments

3 **A.1 Real-world prompt demonstration**

4 Here we describe how we collected and processed a visual, human demonstration in the real-world
5 to treat as a prompt for the trained TTP policy (Fig. 2). Essentially, we collect demonstration
6 pointcloud sequences and manually segment them into different pick-place segments, followed by
7 extracting object states. At each high-level step, we measure the state using three RealSense RGBD
8 cameras[1], which are calibrated to the robot frame of reference using ARTags [2]. The camera
9 output, extrinsics, and intrinsics are combined using Open3D [3] to generate a combined pointcloud.
10 This pointcloud is segmented and clustered to give objects' pose and category using the algorithm
11 from [4] and DBScan. For each object point cloud cluster, we identify the object pose based on the
12 mean of the point cloud. For category information we use median RGB value of the pointcloud,
13 and map it to apriori known set of objects. In the future this can be replaced by more advanced
14 techniques like MaskRCNN [5]. Placement poses are approximated as a fixed, known location, as
15 the place action on hardware is a fixed 'drop' position and orientation. The per step state of the
16 objects is used to create the input prompt tokens used to condition the policy rollout in the real-
17 world, as described in Section 3.2.

18 **A.2 Hardware policy rollout**

19 We zero-shot transfer our policy π trained in simulation to robotic hardware, by assuming low-
20 level controllers. We use a Franka Panda equipped with a Robotiq 2F-85 gripper, controlled using

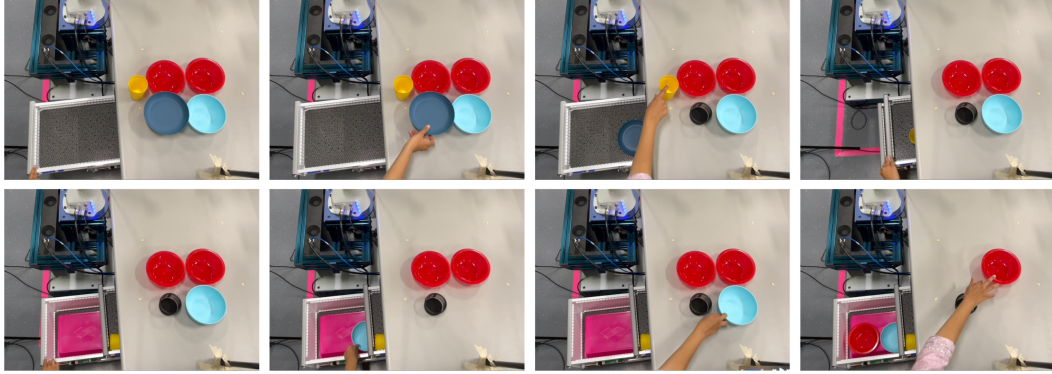


Figure 2: Human demonstration of real-world rearrangement of household dishes.

21 the Polymetis control framework [6]. Our hardware setup mirrors our simulation, with different
 22 categories of dishware (bowls, cups, plates) on a table, a “dishwasher” (cabinet with two drawers).
 23 The objective is to select an object to pick and place it into a drawer (rack) (see Fig. 2).

24 Once we collect the human prompt demonstration tokens, we can use them to condition the learned
 25 policy π from simulation. Converting the hardware state to tokens input to π follows the same
 26 pipeline as the ones used for collecting human demonstrations. At each step, the scene is captured
 27 using 3 Realsense cameras, and the combined pointcloud is segmented and clustered to get object
 28 poses and categories. This information along with the timestep is used to generate instance tokens
 29 as described in Section 2 for all objects visible to the cameras. For visible already placed objects,
 30 the place pose is approximated as a fixed location. The policy π , conditioned on the human demo,
 31 reasons about the state of the environment, and chooses which object to pick. Next, we use a grasp
 32 generator from [7] that operates on point clouds to generate candidate grasp locations on the chosen
 33 object. We filter out grasp locations that are kinematically not reachable by the robot, as well as
 34 grasp locations located on points that intersect with other objects in the scene. Next, we select the
 35 top 5 most confident grasps, as estimated by the grasp generator, and choose the most top-down
 36 grasp. We design an pre-grasp approach pose for the robot which is the same final orientation as
 37 the grasp, located higher on the grasping plane. The robot moves to the approach pose following a
 38 minimum-jerk trajectory, and then follows a straight line path along the approach axes to grasp the
 39 object. Once grasped, the object is moved to the pre-defined place pose and dropped in a drawer.
 40 The primitives for opening and closing the drawers are manually designed on hardware.

41 The learned policy, conditioned on prompt demonstrations, is applied to two variations of the same
 42 scene, and the predicted pick actions are executed. Fig.3 shows the captured image from one of
 43 the three cameras, the merged point cloud and the chosen object to pick and selected grasp for
 44 the same. The policy was successful once with 100% success rate, and once with 75%, shown in
 45 Fig.???. The failure case was caused due to a perception error – a bowl was classified as a plate. This
 46 demonstrates that our approach (TTP) can be trained in simulation and applied directly to hardware.
 47 The policy is robust to minor hardware errors like a failed grasp; it just measures the new state of the
 48 environment and chooses the next object to grasp. For example, if the robot fails to grasp a bowl,
 49 and slightly shifts the bowl, the cameras measure the new pose of the bowl, which is sent to the
 50 policy. However, TTP relies on accurate perception of the state. If an object is incorrectly classified,
 51 the policy might choose to pick the wrong object, deviating from the demonstration preference. In
 52 the future, we would like to further evaluate our approach on more diverse real-world settings and
 53 measure its sensitivity to the different hardware components, informing future choices for learning
 54 robust policies.

55 A.3 Transforming hardware to simulation data distribution

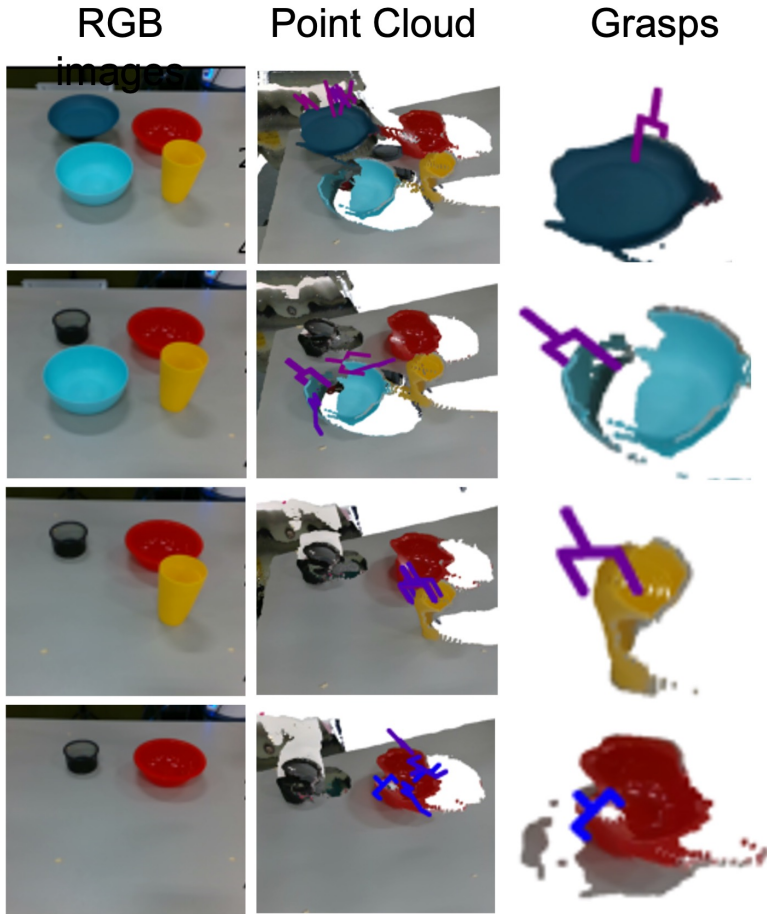


Figure 3: Point cloud and grasps for different objects during policy rollout.

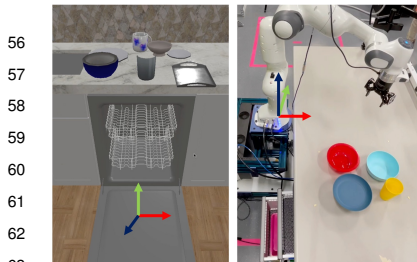


Figure 4: Coordinate Frame of reference in simulation (left) and real world setting (right). Red is x-axis, green is y-axis and blue is z-axis.

the real world coordinates to simulation frame coordinates for running the trained TTP policy on a Franka arm.

We use the semantic work area in simulation and hardware to transform the hardware position coordinates to simulation position coordinates. We measure the extremes of the real workspace by manually moving the robot to record positions and orientations that define the extents of the workspace for the table. The extents of the drawers are measured by placing ARTag markers. We build 3 real-to-sim transformations using the extents for counter, top rack and bottom rack: Let $X \in \mathbb{R}^{3 \times N}$

The policy trained in simulation applies zero-shot to real-world scenarios, but it requires a coordinate transform. Fig. 4 shows the coordinate frame of reference in simulation and real world setting. Since our instance embedding uses the poses of objects, it is dependant on the coordinate frame that the training data was collected in. Since hardware and simulation are significantly different, this coordinate frame is not the same between sim and real. We build a transformation that converts hardware measured poses to the simulation frame of reference, which is then used to create the instance tokens. This ensures that there is no sim-to-real gap in object positions, reducing the challenges involved in applying such a simulation trained policy to hardware. In this section we describe how we convert

76 contain homogeneous xz - coordinates of a work area, along its column, as follows:

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots \\ z^{(1)} & z^{(2)} & \dots \\ 1 & 1 & \dots \end{bmatrix} = [\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \dots] \quad (1)$$

77 As the required transformation from real to simulation involves scaling and translation only, we
 78 have 4 unknowns, namely, $\mathbf{a} = [\alpha_x, \alpha_y, x_{trans}, z_{trans}]$. Here α_x, α_z are scaling factors and
 79 x_{trans}, z_{trans} are translation offset for x and z axis respectively. To solve $X_{sim} = AX_{hw}$, we

80 need to find the transformation matrix $A = \hat{\mathbf{a}} = \begin{bmatrix} \alpha_x & 0 & x_{trans} \\ 0 & \alpha_z & z_{trans} \\ 0 & 0 & 1 \end{bmatrix}$.

$$X_{sim} = \hat{\mathbf{a}}X_{hw} \quad (2)$$

Rewriting the system of linear equations, (3)

$$\Rightarrow \begin{bmatrix} x_{sim}^{(1)} \\ z_{sim}^{(1)} \\ x_{sim}^{(2)} \\ z_{sim}^{(2)} \\ \vdots \end{bmatrix} = \begin{bmatrix} x_{hw}^{(1)} & 0 & 1 & 0 \\ 0 & z_{hw}^{(1)} & 0 & 1 \\ x_{hw}^{(2)} & 0 & 1 & 0 \\ 0 & z_{hw}^{(2)} & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \mathbf{a}^T \quad (4)$$

(5)

81 Let the above equation be expressed as $Y_{sim} = Z_{hw}a^T$ where $Y_{sim} \in \mathbb{R}^{2N \times 1}$, $Z_{hw} \in \mathbb{R}^{2N \times 4}$, and
 82 $a^T \in \mathbb{R}^{4 \times 1}$. Assuming we have sufficient number of pairs of corresponding points in simulation
 83 and real world, we can solve for \mathbf{a} by least squares $\mathbf{a} = (Z_{hw}^T Z_{hw})^{-1} Z_{hw}^T Y_{sim}$. The height y_{sim} is
 84 chosen from a look-up table based on y_{hw} . Once we compute the transformation A , we store it for
 85 later to process arbitrary coordinates from real to sim, as shown below.

```
def get_simulation_coordinates(xyz_hw: List[float], A: np.array) -> List:
    xz_hw = [xyz_hw[0], xyz_hw[2]]
    X_hw = get_homogenous_coordinates(xz_hw)
    X_sim_homo = np.matmul(A, X_hw)
    y_sim = process_height(xyz_hw[1])
    X_sim = [X_sim_homo[0]/X_sim_homo[2], y_sim, X_sim_homo[1]/X_sim_homo[2]]
    return X_sim
```

86 The objects used in simulation training are different from hardware objects, even though they belong
 87 to the same categories. For example, while both sim and real have a small plate, the sizes of these
 88 plates are different. We can estimate the size of the objects based on actual bounding box from the
 89 segmentation pipeline. However, it is significantly out-of-distribution from the training data, due to
 90 object mismatch. So, we map each detected object to the nearest matching object in simulation and
 91 use the simulation size as the input to the policy. This is non-ideal, as the placing might differ for
 92 sim versus real objects. In the future, we would like to train with rich variations of object bounding
 93 box size in simulation so that the policy can generalize to unseen object shapes in the real world.

94 B Simulation Setup

95 B.1 Dataset

96 “Replica Synthetic Apartment 0 Kitchen” consists of a fully-interactive dishwasher with a door and
 97 two sliding racks, an adjacent counter with a sink, and a “stage” with walls, floors, and ceiling.
 98 We use selected objects from the ReplicaCAD [8] dataset, including seven types of dishes (cups,
 99 glasses, trays, small bowls, big bowls, small plates, big plates) which are loaded into the dishwasher.
 100 Fig. 5 shows a human demonstration recorded in simulation by pointing and clicking on desired

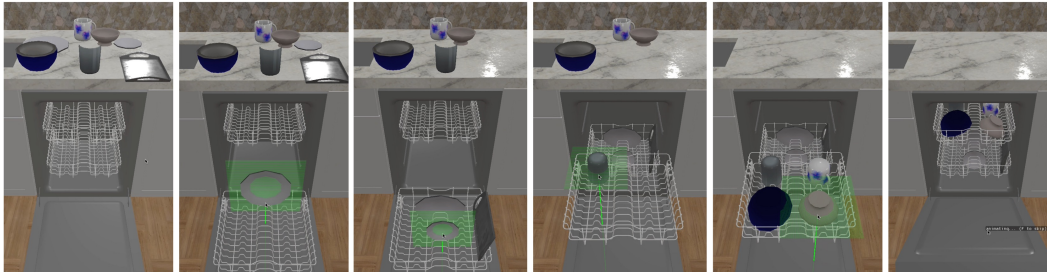


Figure 5: Human demonstration with point and click in simulation

101 object to pick and place. We initialize every scene with an empty dishwasher and random objects
 102 placed on the counter. Next, we generate dishwasher loading demonstrations, adhering to a given
 103 preference, using an expert-designed data generation script. Expert actions include, opening/closing
 104 dishwasher/racks, picking/placing objects in feasible locations or the sink if there are no feasible
 105 locations left. Experts differ in their preferences, and might choose different object arrangements in
 106 the dishwasher.

107 B.2 Expert Preferences

108 We define a preference in terms of expert demonstration ‘prop-
 109 erties’, like which rack is loaded first with what objects? There
 110 are combinatorially many preferences possible, depending on
 111 how many objects we use in the training set. For example, Ta-
 112 ble 1 describes the preferences of dishwasher loading in terms
 113 of three properties - first loaded tray, objects in top and bottom
 114 tray. Each preference specifies properties such as which rack
 115 to load first and their contents. In Table 1, Preferences 1 &
 116 2 vary in the order of which rack is loaded first, while 2 & 3
 117 both load the bottom rack first with similar categories on top
 118 and bottom but with different orderings for these categories.
 119 Other preferences can have different combinations of objects
 120 loaded per rack.

121 To describe a preference, let there be k properties, where each
 122 can take m_k values respectively. For example, a property to
 123 describe preference can be which rack is loaded first, and this
 124 can take two values; either top or bottom rack. The total num-
 125 ber of possible preferences is $G = \prod_{i=1}^k m_i$.

126 In our demonstration dataset, we have 100 unique sessions per preference. Each session can act as
 127 a prompt to indicate preference as well as provide situation for the policy. Each session is about
 128 ~ 30 steps long. With 7 preferences, this leads to $70,000 \times 30 = 2,100,000 \sim 2$ million total
 129 training samples, creating a relatively large training dataset from only 100 unique demonstrations
 130 per preference.

131 Individual task preferences differ in the sequence of expert actions, but collectively, preferences
 132 share the underlying task semantics. For example, the user always opens the dishwasher rack be-
 133 fore loading it for all preferences. By jointly learning over all preferences, our policy can benefit
 134 from cross-preference data to learn task structure, and sparse per-preference data to learn expert
 135 preference.

Table 1: Three example preferences for dishwasher loading. Rack order and their respective contents (ordered by preference).

<i>First?</i>	<i>Top</i>	<i>Bottom</i>
Top	1. cups 2. glasses 3. small bowl	1. big plates 2. small plates 3. trays 4. big bowls
	1. cups 2. glasses 3. small bowl	1. big plates 2. small plates 3. trays 4. big bowl
Bottom	1. small plate 2. glasses 3. cups	1. big bowls 2. trays 3. big plates 4. small bowl

136 B.3 Dynamically appearing objects

137 To add additional complexity to our simulation environment, we simulate a setting with dynami-
138 cally appearing objects later in the episode. During each session, the scene is initialized with $p\%$
139 of maximum objects allowed. The policy/expert starts filling a dishwasher using these initialized
140 objects. After all the initial objects are loaded and both racks are closed, new objects are initialized
141 one-per-timestep to the policy. The goal is to simulate an environment where the policy does not
142 have perfect knowledge of the scene, and needs to reactively reason about new information. The
143 policy reasons on both object configurations in the racks, and the new object type to decide whether
144 to ‘open a rack and place the utensil’ or ‘drop the object in the sink’.

145 C Training

146 In this Section we describe details of the different components of our learning pipeline.

147 C.1 Baseline: GNN

148 **Architecture** We use GNN with attention. The input consists of 12 dimensional attribute inputs
149 (1D-timestep, 3D-category bounding box extents, 7D-pose, 1D-is object or not?) and 12 dimen-
150 sional one-hot encoding for the preference.

```
input_dim: 24  
hidden_dim: 128  
epochs: 200  
batch_size: 32
```

151 **Optimizer** : Adam with $lr = 0.01$ and $weight_decay = 1e - 3$.

152 **Reward function for GNN-RL** Reward function for the RL policy is defined in terms of prefer-
153 ence. The policy gets a reward of +1 every time it predicts the instance to pick that has the category
154 according to the preference order and whether it is placed on the preferred rack.

155 C.2 Our proposed approach: TTP

156 **Architecture** We use a 2-layer 2-head Transformer network for encoder and decoder. The input
157 dimension of instance embedding is 256 and the hidden layer dimension is 512. The attributes
158 contribute to the instance embedding as follows:

```
C_embed: 16  
category_embed_size: 64  
pose_embed_size: 128  
temporal_embed_size: 32  
marker_embed_size: 32
```

159 For the slot attention layer at the head of Transformer encoder, we use:

```
num_slots: 50  
slot_iters: 3
```

160 **Optimizer** We use a batch-size of 64 sequences. Within each batch, we use pad the inputs with 0
161 upto the max sequence length. Our optimizer of choice is SGD with momentum 0.9, weight decay
162 0.0001 and dampening 0.1. The initial learning rate is 0.01, with exponential decay of 0.9995 per
163 10 gradient updates. We used early stopping with patience 100.

164 C.3 Metrics

165 In Section 3 , we presented packing efficiency (PE) and edit distance
 166 (ED) metrics collected on a policy rollout. We present additional metrics
 167 about training progress and rollout here.

168 **Category-token Accuracy** indicates how well the policy can mimic the
 169 expert’s action, given the current state. We monitor training progress by
 170 matching the predicted instance to the target chosen in demonstration
 171 (Fig. 6). We see that TTP is able to predict the same category object to
 172 pick perfectly (accuracy close to 1.0). However, this is a simpler setting
 173 that sequential decision making. During rollout, any error in a state
 174 could create a setting that is out-of-distribution for the policy. Thus,
 175 category token accuracy sets an upper bound for rollout performance,
 176 that is, while having high category token accuracy is necessary, it is not
 177 sufficient for high packing efficiency and inverse edit distance.

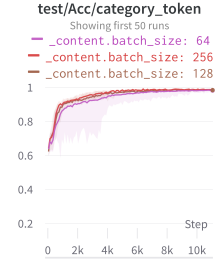
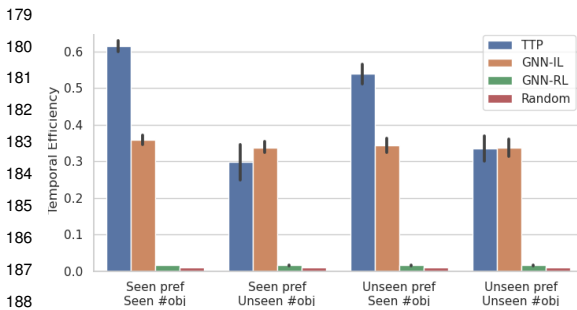


Figure 6: Category level accuracy grouped by batch size for prompt-situation training.

178 **Temporal efficiency:**



189 Figure 7: TE (SPL) metric for held-out test settings.

190
 191 between 0 to 1, and higher is better. This value will be equal to or lower than the packing efficiency.
 192 This especially penalizes policies that present a ‘looping’ behavior, such as repeatedly open/close
 193 dishwasher trays, over policies that reach a low PE in shorter episodes (for example, by placing
 194 most objects in the sink). Fig 7 shows the temporal efficiency or SPL over our 4 main held-out test
 195 settings.

196 **D Additional Ablation Experiments**

197 In Section ?? we presented ablation experiments over number of demonstrations per preference
 198 used for training, and the number of unique preferences used. In this Section, we present additional
 199 ablation experiments over the design of instance encodings in TTP. Additionally, we also present
 200 results where we increase the temporal context of TTP and study its effect on performance.

201 **D.1 Design of Instance Encoding**

202 **How much does temporal encoding design matter?** Fig. 8a shows that learning an embedding
 203 per timestep or expanding it as fourier transformed vector of sufficient size achieves high success.
 204 On the other hand, having no timestep input shows slightly lower performance. Timestep helps in
 205 encoding the order of the prompt states. The notion of timestep is also incorporated by autoregres-
 206 sive masking in both the encoder and the decoder.

207 **How much does category encoding design matter?** In our work, we represent category as the
 208 extents of an objects’ bounding box. An alternative would be to denote the category as a discrete
 209 set of categorical labels. Intuitively, bounding box extents captures shape similarity between objects
 210 and their placement implicitly, which discrete category labels do not. Fig. 8b shows that fourier
 211 transform of the bounding box achieves better performance than discrete labels, which exceeds the
 212 performance with no category input.

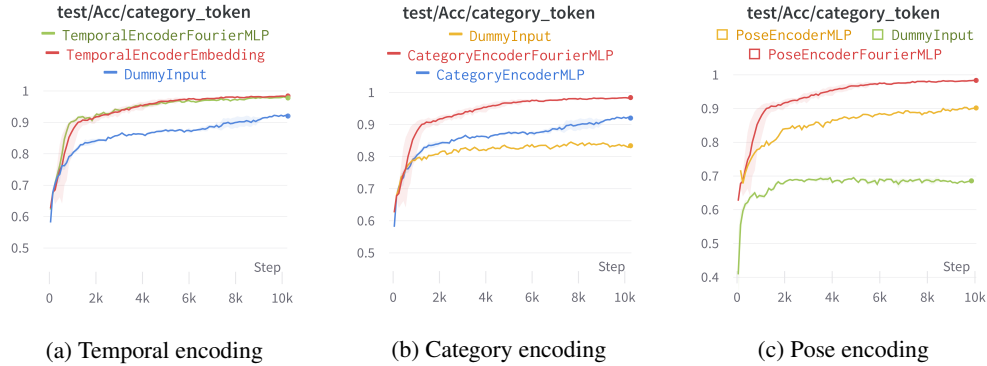


Figure 8: [Left-to-Right] Comparing different design choices of attribute encoders in terms of category token accuracy on held-out test prompt-situation session pairs.

213 **How much does pose encoding design matter?** We encode pose as a 7-dim vector that includes
 214 3d position and 4d quaternion. Fig. 8c shows that the fourier transform of the pose encoding
 215 performs better than feeding the 7 dim through MLP. Fourier transform of the pose performs better
 216 because such a vector encodes the fine and coarse nuances appropriately, which otherwise either
 217 require careful scaling or can be lost during SGD training.

218 D.2 Markov assumption on the current state in partial visibility scenarios

219 Dynamic settings, as used in our simulation, can be
 220 partially observable. For example, when the rack is
 221 closed, the policy doesn't know whether it is full or
 222 not from just the current state. If a new object ar-
 223 rives, the policy needs to decide between opening
 224 the rack if there is space, or dropping the object in
 225 sink if the rack is full. In such partially observed set-
 226 tings, the current state may or may not contain all the
 227 information needed to reason about the next action.
 228 However, given information from states in previous
 229 timesteps, the policy can decide what action to take
 230 (whether to open the rack or directly place the ob-
 231 ject in the sink). With this in mind, we train a single
 232 preference pick only policy for different context his-
 233 tory. As shown in Fig. 10, context window of size k
 234 processes the current state as well as k predecessor
 235 states, that is, in total $k + 1$ states.

236 Let context history k refer to the number of previous
 237 states included in the input. Then the input is a se-
 238 quence of previous k states' instances (including the
 239 current state), as shown in Fig. 10.

240 Fig 9 shows that TTP gets $> 90\%$ category level prediction accuracy in validation for all context
 241 windows. While larger context windows result in faster learning at the start of the training, the
 242 asymptotic performance of all contexts is the same. This points to the dataset being largely visible,
 243 and a single context window capturing the required information. In the future, we would like to
 244 experiment with more complex settings like mobile robots, which might require a longer context.

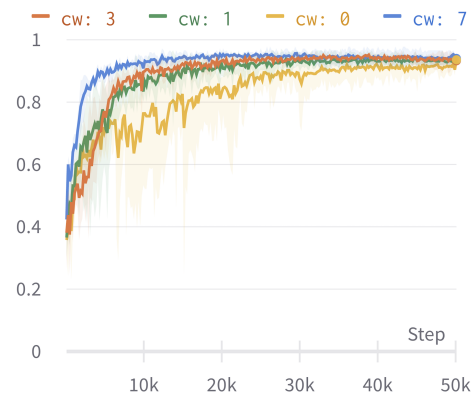


Figure 9: Plot showing category level accuracy for the held-out test sessions for single preference training with context windows. While larger context window size learns faster, the asymptotic performance for all context windows converges in our setting.

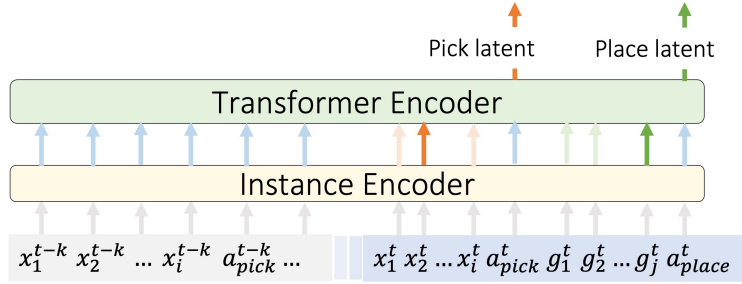


Figure 10: Processing with previous Context History k

245 E Limitations and Future scope

246 In Section 6, we briefly discussed the limitations and risks. Here we enlist more details and highlight
 247 future directions.

248 **Pick grasping depends on accurate segmentation and edge detection** Grasping policy depends
 249 on quality of segmentation and edge detection of the selected object. Due to noise in calibration,
 250 shadows and reflections, there are errors in detecting the correct edge to successfully grasp the
 251 object. For example, it is hard to grasp a plate in real setting. Plate is very close to the ground and
 252 the depth cameras cannot detect a clean edge for grasping. Therefore, in our work, we place the
 253 plate on an elevated stand for easy grasping. Grasping success also depends on the size and kind of
 254 gripper used.

255 **Placement in real setting** For placement, the orientation of final pose is often different from
 256 initial pose and may require re-grasping. The placement pose at final settlement is different from the
 257 robot’s end-effector pose while releasing the object from its grasp. Similar to picking, placement
 258 accuracy will largely depend on appropriate size and shape of gripper used. Due to these reasons,
 259 placement in real world is an open challenging problem and we hope to address this future work.

260 **Hardware pipeline issues due to calibration** The resulting point cloud is generated noisy due to
 261 two reasons. First, incorrect depth estimation due to camera hardware, lighting conditions, shadows
 262 and reflections. Second, any small movements among cameras that affects calibration. If we have
 263 a noisy point cloud, it is more likely to have errors in subsequent segmentation and edge detection
 264 for grasp policy. Having sufficient coverage of the workspace with cameras is important to mitigate
 265 issues due to occlusions and incomplete point clouds.

266 **Incomplete information in prompt** The prompt session may not contain all the information to
 267 execute on the situation. For example, in a prompt session there might be no large plates seen, which
 268 is incomplete/ambiguous information for the policy. This can be mitigated by ensuring complete
 269 information in prompt demo or having multiple prompts in slightly different initialization.

270 **References**

- 271 [1] L. Keselman, J. Iselin Woodfill, A. Grunnet-Jepsen, and A. Bhowmik. Intel realsense stereo-
272 scopic depth cameras. In *Proceedings of the IEEE conference on computer vision and pattern*
273 *recognition workshops*, pages 1–10, 2017.
- 274 [2] M. Fiala. Artag, a fiducial marker system using digital techniques. In *2005 IEEE Computer*
275 *Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 2, pages
276 590–596. IEEE, 2005.
- 277 [3] Q.-Y. Zhou, J. Park, and V. Koltun. Open3d: A modern library for 3d data processing. *arXiv*
278 *preprint arXiv:1801.09847*, 2018.
- 279 [4] Y. Xiang, C. Xie, A. Mousavian, and D. Fox. Learning rgb-d feature embeddings for unseen
280 object instance segmentation. *arXiv preprint arXiv:2007.15157*, 2020.
- 281 [5] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE*
282 *international conference on computer vision*, pages 2961–2969, 2017.
- 283 [6] Y. Lin, A. S. Wang, G. Sutanto, A. Rai, and F. Meier. Polymetis. [https://](https://facebookresearch.github.io/fairo/polymetis/)
284 facebookresearch.github.io/fairo/polymetis/, 2021.
- 285 [7] H.-S. Fang, C. Wang, M. Gou, and C. Lu. Graspnet-1billion: A large-scale benchmark for
286 general object grasping. In *Proceedings of the IEEE/CVF conference on computer vision and*
287 *pattern recognition*, pages 11444–11453, 2020.
- 288 [8] A. Szot, A. Clegg, E. Undersander, E. Wijmans, Y. Zhao, J. Turner, N. Maestre, M. Mukadam,
289 D. Chaplot, O. Maksymets, A. Gokaslan, V. Vondrus, S. Dharur, F. Meier, W. Galuba, A. Chang,
290 Z. Kira, V. Koltun, J. Malik, M. Savva, and D. Batra. Habitat 2.0: Training home assistants to
291 rearrange their habitat. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- 292 [9] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik,
293 R. Mottaghi, M. Savva, et al. On evaluation of embodied navigation agents. *arXiv preprint*
294 *arXiv:1807.06757*, 2018.