

1 We thank the reviewers for their insightful comments. We are glad reviewers acknowledge our promising approach of  
2 using hierarchical control and relational information to generalize bimanual manipulation imitation learning **R1, R2,**  
3 **R3, R4.** We are encouraged by comments about the soundness of our experiments and theory **R1, R3.** We will fix the  
4 minor issues and typos in the updated version. Below we clarify our approach and address specific concerns.

5 **R1, R2 *Related work and baselines*** To the best of our knowledge, using deep imitation learning for bimanual  
6 manipulation have only involved MLP [1] and RNN [2]. [4] uses graph representation for visual imitation learning.  
7 MLP did not suit our task as we wanted to predict trajectory sequences. The GRU-GRU model presented in Table 1  
8 was an implementation of RNN, which we adapted to our simulation environment. Graph neural networks with skip  
9 connections in the context of bimanual manipulation have not been previously investigated. We also experimented with  
10 GAIL [3] but it did not perform well. Errors for the first primitive were at least 6 times higher than HDR-IL. This is  
11 likely due to the low variances of our expert demonstrations, making the policy distribution difficult to fit.

12 **R1, R4 *29% accuracy rather low...*** The success rate is low because the task is intended to test the limits. In our  
13 high-precision table lifting task, we showed 100% success rate in Table 1. This demonstrates our HDR-IL model greatly  
14 improved generalization compared to prior methods. To deepen our investigation, we designed the more challenging  
15 peg-in-hole task to test multiple generalizations in far out (up to 130) time steps.

16 **R1 *...an error of 5cm with HDR-IL seems unreasonably high.*** The Euclidean distance errors shown in Table 1 **are**  
17 **normalized** to per datapoint position. These errors do not necessarily reflect the simulation task performance. Task  
18 success is largely determined by the accuracy at the end of a primitive. Errors in our model tend to be at the beginning  
19 as suggested by the trajectories in Figure 4. Furthermore, Table 2 of the Appendix show the 5cm error is largely driven  
20 by the Extend primitive, which is less crucial for the task success. The Extend primitive has the largest step sizes  
21 between datapoints, as illustrated in Figure 3(a). The bigger steps leads to bigger prediction errors.

22 **R1 *...actual pose of one arm...relative pose of the second arm...*** We ran tests and found using relational data does not  
23 generalize as well as graph structure with absolute poses. Euclidean distances were  $8.15 \pm 4.71$  with success rate 16%  
24 compared to  $5.64 \pm 5.17$  and 72% for ResInt in Table 1. We will add these to our baseline in an updated version.

25 **R1 *...effect of orientation...*** Orientation of the grippers were included in the prediction to improve inverse kinematics  
26 accuracy. We will quantify orientation errors of the gripper quaternions in an updated version. The generalization of the  
27 model to different table starting orientations is a topic for future study.

28 **R1 *Not clear if...specific to bimanual manipulation*** Our method is NOT specific to bimanual manipulation, but the  
29 complex dependencies in the bimanual manipulation setting highlight the value of graphs in our experiments.

30 **R2 *...why did the ground truth demonstration fail...*** The failures in demonstration are due to approximations of the  
31 inverse kinematics (IK) solver. Such failure cases are rare ( $<0.7\%$  of 2500 demonstrations). The IK solver sometimes  
32 outputs unusual movements between poses which cause the table to be dropped. The arm trajectory in the demonstration  
33 will change accordingly. Our model can account for these uncertainties in demonstrations because of the stochastic  
34 sampling design shown in Figure 2.

35 **R2 *...limited to the provision of the graph structure.*** Providing the graph structure is a common way to encode inductive  
36 bias e.g. [4], not a limiting factor. Learning the graph structure is related to model selection and causal reasoning, which  
37 is out of scope of this paper. However, we did experiment with a few different graph structures before arriving at our  
38 fully connected graph with attention.

39 **R2 *...does not show off how complex two handed settings...*** Our simulations include complex **physical contacts and**  
40 **friction.** We did not use a flexible table. The peg-in-hole task was specifically designed to increase the task complexity:  
41 assembling the table before lifting requires more coordination between the two arms. We showed a failed attempt at  
42 lifting the table using one arm in our supplementary videos. Doing the same for either piece in the peg-in-hole task will  
43 have the same result. A hammer and nail task would be an even more complicated experiment for future work.

44 **R2 *...manually crafted primitives...contradicting to learning...*** This is NOT a contradiction. The primitives are  
45 constructed based on the control commands executed in simulation. During training, the primitives are labeled. Our  
46 model learns to compose the primitives for high level planning. It also learns the low level dynamics for each primitive.  
47 During inference, the model predicts the primitives based on the demonstration trajectories.

48 [1] R. Chitnis, et al., “Efficient Bimanual Manipulation Using Learned Task Schemas,” *ICRA*, 2020

49 [2] D. Rakita, et al., “Shared Control–Based Bimanual Robot Manipulation,” *Science Robotics*, 2019

50 [3] J. Ho, et al., “Generative Adversarial Imitation Learning,” *NeurIPS*, 2016

51 [4] M. Sieb, et al., “Graph-structured Visual Imitation,” *CoRL*, 2020