

---

# Hierarchical Abstraction for Combinatorial Generalization in Object Rearrangement

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Object rearrangement is a challenge for embodied agents because solving these  
2        tasks requires generalizing across a combinatorially large set of underlying entities  
3        that take the value of object states. Worse, these entities are often unknown and  
4        must be inferred from sensory percepts. We present a hierarchical abstraction  
5        approach to uncover these underlying entities and achieve combinatorial general-  
6        ization from unstructured inputs. By constructing a factorized transition graph  
7        over clusters of object representations inferred from pixels, we show how to learn  
8        a correspondence between intervening on states of entities in the agent’s model  
9        and acting on objects in the environment. We use this correspondence to develop  
10       a method for control that generalizes to different numbers and configurations of  
11       objects, which outperforms current offline deep RL methods when evaluated on a  
12       set of simulated rearrangement and stacking tasks.

## 13    1 Introduction

14    A core property of intelligence is the ability to re-purpose previously acquired knowledge for solving  
15    new problems, but how to identify symmetries for factorizing modeling and control into independent  
16    components has been a challenge for AI systems, which either assume human-defined abstractions to  
17    begin with [11, 13] or learn monolithic representations with no clear mechanism for reuse [19, 30].

18    The problem of *object rearrangement* offers an intuitive setting for studying this problem of knowl-  
19    edge reuse. Because the space of object configurations is combinatorially large, solving novel  
20    rearrangement problems requires the agent to recognize that the same action for moving an object  
21    from one location to another can be reused for different objects in different contexts. We can formulate  
22    this problem as an offline goal-conditioned reinforcement learning (RL) problem, where the agent is  
23    trained on a dataset of sensorimotor interactions and is evaluated on rearranging objects to satisfy  
24    constraints depicted in a goal image. The research questions are, therefore, how can we enable the  
25    agent to represent abstractions of objects in a way that is amenable for such reuse in control, and how  
26    can we enable it to discover these abstractions on its own directly from the sensorimotor interface?

27    Our main contribution is a method called Hierarchical Abstraction (HA) that offers an answer to the  
28    above questions. The key idea is to represent objects as independent and symmetric latents, which we  
29    call entities, and factorize each entity  $h^k$  into two parts, its identity  $z^k$  and its state  $s^k$ . The identity  
30    is the part that is not affected by action and the state is the part that is. The modeling component  
31    of our approach abstracts the dataset of sensorimotor interactions into a graph over state transitions  
32    of individual entities. By construction, a state transition is agnostic to the identity of the entity or  
33    the context entities present in the sensorimotor interaction from which it was learned, enabling it  
34    to be reused across multiple entities and contexts. The control component of our approach then  
35    re-composes sequences of state transitions to solve new rearrangement problems by inferring what  
36    state transition can be taken given only the current and goal image observations.

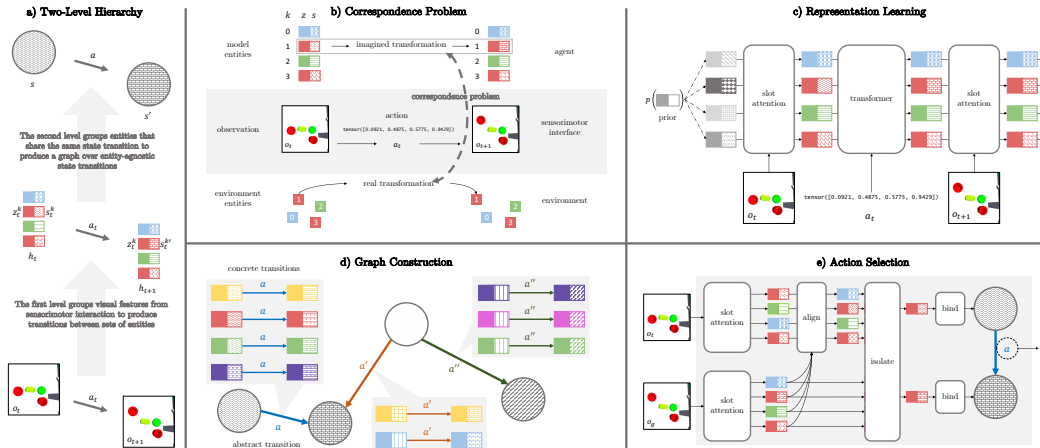


Figure 1: **Overview** (a) HA abstracts video interactions into a graph over state transitions of individual entities with a two-level hierarchy. (b) To do so requires learning representations of entities in a way that exhibits a correspondence between how model entities and how real objects transform under action. (c) The first level of the hierarchy abstracts visual features into transitions over sets of entities. (d) The second level abstracts these transitions into a graph of transitions over individual states. (e) HA decomposes the rearrangement task into a subtask per entity and chooses actions by returning the action tagged to the edge between nodes in the graph.

37 In contrast to other works in object-centric learning [14–16, 21, 24, 31, 34, 40] which evaluate the  
 38 quality of inferred entities via segmentation metrics, we evaluate our inferred entities by *how well*  
 39 *they can reused for solving tasks*, which more directly assesses what we want an agent’s object  
 40 representations to do. Kulkarni et al. [23], Veerapaneni et al. [36] also evaluate on control tasks.  
 41 Whereas Kulkarni et al. [23] considers how object representations improve exploration, we consider  
 42 the offline setting which requires zero-shot generalization. The shooting-based planning method  
 43 in Veerapaneni et al. [36] suffers from compounding errors as other learned single-step models  
 44 do [20], but our non-parametric approach of composing previously seen transitions enables us to plan  
 45 for longer horizons. Indeed, our results show that HA outperforms both state-of-the-art offline RL  
 46 methods and slot-based planning methods across three simulated object rearrangement problems.

## 47 2 Hierarchical Abstraction

48 We present our method for solving object rearrangement problems from pixels. Given a transition  
 49  $o, a \rightarrow o'$  from the training set in which a single object  $k$  has been moved, the modeling problem is  
 50 to represent the state transition  $s^k \rightarrow s^{k'}$  as decoupled from the identity  $z^k$  of the affected entity and  
 51 as decoupled from the other context entities  $h^{\neq k}$  in the scene. The control problem is to compose  
 52 previously seen state transitions to satisfy goal constraints on new object configurations. Though we  
 53 do not assume the number of action primitives is finite, we assume the number of observed locations  
 54 can be clustered into a set of finite groups.

### 55 2.1 Modeling

56 Our modeling approach abstracts the training dataset of action-conditioned videos into a factorized  
 57 graph over state transitions that can be reused across different rearrangement problems. Constructing  
 58 such a graph decomposes into two representation learning problems: the first concerns how to  
 59 represent objects in a way that exposes their commonalities and the second concerns how to represent  
 60 actions in a way that enables them to be reused across these commonalities. Our key insight is that  
 61 these two problems are both different instances of the same fundamental problem of partitioning a  
 62 space into discrete groups, but at different levels of abstraction: learning how to cluster visual features  
 63 into entities (Fig. 1b) and learning how to cluster transitions over entity-sets into transitions over  
 64 individual states (Fig. 1c). Our model thus implements a two-level abstraction hierarchy (Fig. 1a).

65 **Level 1: abstracting visual features into sets of entities** The goal of the first level is to map  
 66 a transition over two video frames  $o_1, a_1 \rightarrow o_2$  into a transition over entity-sets  $(h_1^1, \dots, h_1^K) \rightarrow$   
 67  $(h_2^1, \dots, h_2^K)$  (Fig. 1c). The criteria our entity-set transitions should satisfy are (1) entities represent

68 different objects in the scene, i.e. they are *independent* and *symmetric* and (2) only the state  $s^k$  of  
 69 entity  $k$  changes, not its identity  $z^k$ , i.e. the entity is *factorized*. Because each video transition depicts  
 70 only a single object being moved, we also want (3) the transition to be *sparse*, i.e. only one entity  $h^k$   
 71 changes in the transition and the rest  $h^{\neq k}$  remain unchanged.

72 Our approach manually enforces independence, symmetry, and factorization but lets sparsity emerge.  
 73 Specifically, we formulate the inference of entity-sets as the application of an equivariant action-  
 74 conditioned sequential Bayesian filter with a mixture model as the latent state, where entity representa-  
 75 tions are the parameters of the mixture components. Given the connection between the slots produced  
 76 by slot attention (SA) [25] and mixture components Chang et al. [6], we concretely implement this  
 77 filter by temporally evolving the slots produced by the slot attention module of the state-of-the-art  
 78 SLATE architecture [32] with a transformer decoder (TD) [28, 35]. With  $\mathbf{h} = (h^1, \dots, h^K)$ ,

$$\mathbf{h}_{t+1} = SA(\hat{\mathbf{h}}_t, o_{t+1}) \quad \hat{\mathbf{s}}_{t+1} = TD(\text{queries} = \mathbf{s}_t, \text{keys/values} = [\mathbf{s}_t, a_t]),$$

79 where  $\hat{\mathbf{h}}_{t+1} = [\mathbf{z}_t, \hat{\mathbf{s}}_{t+1}]$  and  $\hat{\mathbf{h}}_0$  are independent samples from a Gaussian. This approach enforces  
 80 independence because  $\hat{\mathbf{h}}_0$  are independent, symmetry because the transformer is equivariant, and  
 81 factorization because  $\hat{\mathbf{z}}_{t+1}$  is a copied from  $\mathbf{z}_t$ . Empirically, we find that this architecture models  
 82 sparse changes in the slots (Fig. 3). We call this implementation **dynamic SLATE** (dSLATE), but  
 83 our contribution of how structurally we enforce the above three criteria remains invariant to *how* the  
 84 sequential Bayesian filter is implemented, which we expect will advance beyond SA and TD.

85 **Level 2: abstracting transitions over sets of entities into individual state transitions** The goal  
 86 of the second level is to construct a state transition graph from a buffer of entity-set transitions  
 87  $\{(\mathbf{h}_1, \dots, \mathbf{h}_T)\}_{n=1}^N$  produced from the first level (Alg. 1, Fig. 1d). Specifically, we seek each edge  
 88 of this graph to represent a state transition that has been shared across multiple entities in different  
 89 contexts in the training dataset, such that we can reuse such state transitions when choosing actions  
 90 for solving new rearrangement problems. Given our assumption that observed locations can be  
 91 grouped into a finite set of clusters, this goal translates into a problem of clustering state transitions,  
 92 and the problem of clustering state transitions reduces to clustering the states before and after the  
 93 transition. We treat each cluster centroid as a node in the state transition graph, and an edge between  
 94 nodes is tagged with the single action that transforms one node’s state to another’s.

95 To create the nodes, the first  
 96 step is to identify the entity  $h^k$   
 97 that dSLATE predicted was af-  
 98 fected by  $a_t$  in each transition  
 99  $(\mathbf{h}_t, a_t, \mathbf{h}_{t+1})$ . Our isolate pro-  
 100 cedure achieve this by solving  
 101  $k = \arg \max_{k' \in \{1, \dots, K\}} d(s_t^{k'}, s_{t+1}^{k'})$   
 102 to identify the index of the entity  
 103 whose state has most changed  
 104 during the transition, where  $d(\cdot, \cdot)$   
 105 is a distance function, detailed  
 106 in the Appendix. Having the  
 107 buffer of transitions over entity  
 108 sets  $\{(\mathbf{h}_1, \dots, \mathbf{h}_T)\}_{n=1}^N$  to a buffer  
 109 of individual entity transitions  
 110  $\{(h_t^k, a_t, h_{t+1}^k)\}_{n=1}^N$ , and thus a  
 111 buffer of individual state transitions  
 112  $\{(s_t^k, a_t, s_{t+1}^k)\}_{n=1}^N$ , the second  
 113 step is to cluster similar transitions  
 114 together by clustering the states

---

#### Algorithm 1 Building the Graph

---

```

1: input model, buffer
2: for  $\{(o_t, a_t, o_{t+1})\}_n$  in buffer do
3:   # infer entities from transition
4:    $\{(\mathbf{h}_t, a_t, \mathbf{h}_{t+1})\}_n \leftarrow \text{model}(\{o_t, a_t, o_{t+1}\}_n)$ .
5:   # identify which entity changed in transition
6:    $\{(h_t^k, a_t, h_{t+1}^k)\}_n \leftarrow \text{isolate}(\{(\mathbf{h}_t, a_t, \mathbf{h}_{t+1})\}_n)$ 
7: end for
8: # partition transitions by clustering entities
9:  $\{s_*\}_{m=1}^M \leftarrow \text{cluster}(\{(s_t^k, a_t, s_{t+1}^k)\}_{n=1}^N)$ 
10: # transitions between clusters are edges
11: initialize graph with nodes  $s_*^{[m]}$ , for  $m \in [1 : M]$ 
12: for each  $\{(h_t^k, a_t, h_{t+1}^k)\}_n$  do
13:   # infer cluster assignments
14:    $[i], [j] \leftarrow \text{bind}(h_t^k), \text{bind}(h_{t+1}^k)$ 
15:   # tag edge with action  $a_t$ 
16:   graph.edges $[i, j] \leftarrow \text{create-edge}(s_*^{[i]} \xrightarrow{a_t} s_*^{[j]})$ 
17: end for
18: return graph

```

---

115 before and after each transition. Our cluster procedure achieves this by clustering over all states  
 116 together, which is sufficient since the end state of transition is the starting state of another. A proper  
 117 clustering should return centroids  $\{s_*\}_{m=1}^M$  that all correspond to actual states of entities, which we  
 118 hypothesize will correspond to actual locations in the training set based on how dSLATE is trained.

119 The purpose of edges between the nodes is to record an action that transforms an entity from one  
 120 state to another. For each entity transition  $(h_t^k, a_t, h_{t+1}^k)$  we bind  $h_t^k$  and  $h_{t+1}^k$  to their associated  
 121 nodes  $s_*^{[i]}$  and  $s_*^{[j]}$  and create an edge between  $s_*^{[i]}$  and  $s_*^{[j]}$  tagged with action  $a$ , overwriting previous

edges based on the assumption that with a proper clustering there should only be one action primitive per pair of nodes. Applying the `bind` procedure to an entity  $h^k$  returns the index of the cluster of the centroid  $s_*$  that is nearest to the entity’s state  $s^k$ . For our experiments `cluster` is implemented with K-means and `bind` is implemented with nearest neighbors using the same distance metric used for K-means (see Appendix), but other clustering algorithms and distance metrics can also be used.

## 2.2 Control

Having constructed a graph of state transitions, our approach for control re-composes sequences of state transitions to solve new rearrangement problems. Specifically, the agent decomposes the rearrangement problem into a set of per-entity subproblems, searches the transition graph for a transition that transforms the current entity’s state to its goal state, and executes the action tagged with this transition in the environment. This problem decomposition is possible because the transitions in our graph are constructed to be agnostic to identity and context, enabling different rearrangement problems to share solutions to the same subproblems. The core challenge in deciding which transitions to compose is in determining which transitions are *possible* to compose. That is, the agent must determine which nodes in the graph correspond to the given goal constraints and which nodes correspond to the entities in the current observation, but the current entities  $\mathbf{h}_t$  and goal constraints  $\mathbf{h}_g$  must themselves be inferred from the current and goal observations  $o_t$  and  $o_g$ , requiring the agent to infer both what to do and how to do it purely from its sensorimotor interface.

---

### Algorithm 2 Action Selection

---

```

1: given model, graph
2: input goal  $o_g$ , observation  $o_t$ 
3: # infer goal constraints and current entities
4:  $\mathbf{h}_g, \mathbf{h}_t \leftarrow \text{model}(o_g), \text{model}(o_t)$ 
5: align entity indices of  $\mathbf{h}_t$  with those of  $\mathbf{h}_g$ 
6:  $\pi \leftarrow \text{align}(\mathbf{h}_t, \mathbf{h}_g)$ 
7: permute indices of  $\mathbf{h}_t$  according to  $\pi$ 
8:  $\mathbf{h}_t \leftarrow (h_t^{\pi[1]}, \dots, h_t^{\pi[K]})$ 
9: identify  $k$ th goal constraint to satisfy next
10:  $k \leftarrow \text{select-constraint}(\mathbf{h}_t, \mathbf{h}_g)$ 
11: infer cluster assignments
12:  $[i], [j] \leftarrow \text{bind}(h_t^k), \text{bind}(h_g^k)$ 
13: action that transforms node  $[i]$  to node  $[j]$ 
14: return graph.edges $[i, j]$ .action

```

---

Our approach takes four steps (Alg. 2, Fig. 3). The first step applies dSLATE to infer  $\mathbf{h}_t$  and  $\mathbf{h}_g$  from  $o_t$  and  $o_g$ . The second step permutes the indices of  $\mathbf{h}_t$  and  $\mathbf{h}_g$  to align the identities of  $h_t^k$  with  $h_g^k$ . This step is necessary because the permutation invariance of the entities does not guarantee that the entities with the same identities will have the same index  $k$ , i.e.  $z_t^k$  might not correspond to  $z_g^k$ . We implement this using the Hungarian algorithm that minimizes the matching cost between  $\mathbf{z}^t$  and  $\mathbf{z}^g$ . The third step selects which goal constraint  $h_g^k$  to satisfy next. In the case where objects can move independently, we implement this `select-constraint` procedure by determining which constraint  $h_g^k$  has the highest difference in state with its counterpart  $h_t^k$ , which reduces to solving the same `argxmax` problem as in `isolate`. We discuss how we handle more complex dependencies in the Appendix. Lastly, we `bind`  $h_t^k$  and  $h_g^k$  to the graph and return the action tagged to the edge between their respective nodes. If an edge does not exist between the nodes, we simply take a random action.

## 3 Experiments

The crucial test for knowledge reuse is to solve a set of tasks from a part of the combinatorial space of object rearrangement problems that are disjoint from the part of the space given for training, i.e. generalizing to different numbers of objects. We measure the “fractional success rate:” the mean change in the number of satisfied constraints divided by the number of initially unsatisfied constraints for each episode. In the *complete* setting, all objects have an associated constraint, whereas in the *partial* setting, only a subset of objects have associated constraints, and this is reflected in the goal image (Fig. 2b). We consider three challenging environments: *block-rearrange*, *robogym-rearrange*, and *block-stacking*. We compare against five baselines. We implement a version of MPC [36] by replacing our factorized graph search with CEM. We investigate an ablation of HA that constructs a monolithic graph over transitions between sets of slot states rather than individual slot states (MGS). We also compare with state-of-the-art pixel-based behavior cloning (BC) and implicit Q-learning (IQL) implementations based off of [22]. Our last baseline just takes random actions (Rand).

**Results** Figure 4 shows that HA performs significantly better than the baselines in combinatorial generalization (about a 5-10x improvement). Notably, the MPC baseline performs poorly because its rollouts are poor, and it is significantly more computationally expensive to run (11 hours instead of

175 20 minutes). Fig. 3a shows a t-SNE plot of the factors over 100 batches in the training set, where  
 176 different colors label different clusters, with several clusters annotated with the average slot attention  
 177 mask for the entities in the cluster. The right panel of Fig. 3 walks through the steps our algorithm  
 178 takes to select an action, visualizing the align, isolate, and bind procedures, showing that the  
 179 entities inferred with dSLATE intuitively capture objects that can be independently acted upon.

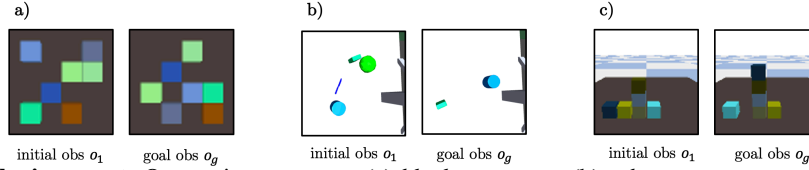


Figure 2: **Environments** Our environments are (a) *block-rearrange*, (b) *robogym-rearrange*, and (c) *block-stacking*. (a) shows the complete specification of goal constraints for all objects; (b) shows a partial specification of goal constraints. (c) shows the case where some goal constraints are already satisfied in the initial observation.

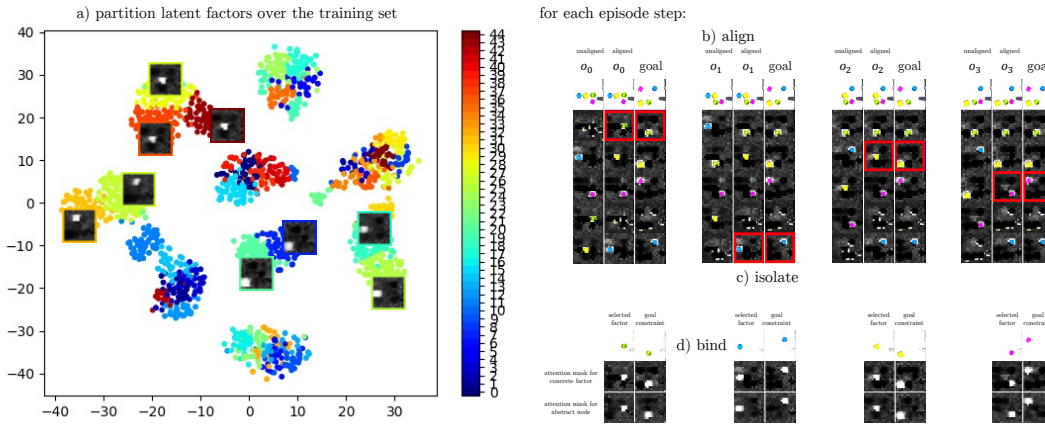


Figure 3: **Qualitative walk-through.** (a) shows the slot attention masks of the centroids from clustering states in *robogym-rearrange* (b) shows a walk-through of how HA selects the next action.

Figure 4: Generalizing from an offline ataset of 4 objects to solving rearrangement tasks with 4-7 objects, evaluated for 100 episodes across 10 seeds, showing standard error.

Table 1: *block-rearrange*, complete specification.

Method	4	5	6	7
Ours	<b>0.94</b> ± 0.01	<b>0.93</b> ± 0.00	<b>0.93</b> ± 0.00	<b>0.89</b> ± 0.00
Rand	0.06 ± 0.02	0.07 ± 0.03	0.07 ± 0.03	0.08 ± 0.03
MPC	0.16 ± 0.06	0.12 ± 0.04	0.11 ± 0.04	0.10 ± 0.03
MGS	0.07 ± 0.03	0.06 ± 0.02	0.07 ± 0.02	0.08 ± 0.03
IQL	0.07 ± 0.01	0.03 ± 0.00	0.02 ± 0.00	0.02 ± 0.00
BC	0.03 ± 0.00	0.02 ± 0.00	0.01 ± 0.00	0.01 ± 0.00

Table 2: *block-stacking*, complete specification.

Method	4	5	6	7
Ours	<b>0.58</b> ± 0.02	<b>0.48</b> ± 0.02	<b>0.35</b> ± 0.01	<b>0.24</b> ± 0.01
Rand	0.02 ± 0.01	0.01 ± 0.01	0.02 ± 0.01	0.02 ± 0.01
MPC	0.03 ± 0.01	0.01 ± 0.01	0.02 ± 0.01	0.01 ± 0.01
MGS	0.13 ± 0.00	0.01 ± 0.01	0.02 ± 0.01	0.01 ± 0.01
IQL	0.19 ± 0.02	0.12 ± 0.01	0.13 ± 0.02	0.07 ± 0.01
BC	0.21 ± 0.01	0.13 ± 0.01	0.13 ± 0.01	0.08 ± 0.01

Table 3: *robogym-rearrange*, complete specification.

Method	4	5	6	7
Ours	<b>0.64</b> ± 0.01	<b>0.47</b> ± 0.01	<b>0.49</b> ± 0.01	<b>0.41</b> ± 0.01
Rand	0.01 ± 0.0	0.01 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
MPC	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
MGS	0.01 ± 0.0	0.01 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
IQL	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
BC	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0

Table 4: *block-stacking*, partial specification.

Method	4	5	6	7
Ours	<b>0.34</b> ± 0.01	<b>0.27</b> ± 0.01	<b>0.21</b> ± 0.01	<b>0.15</b> ± 0.01
Rand	0.03 ± 0.01	0.03 ± 0.01	0.03 ± 0.01	0.03 ± 0.01
MPC	0.04 ± 0.02	0.03 ± 0.01	0.02 ± 0.01	0.04 ± 0.01
MGS	0.03 ± 0.01	0.03 ± 0.01	0.02 ± 0.01	0.02 ± 0.01
IQL	0.14 ± 0.01	0.06 ± 0.01	0.04 ± 0.00	0.03 ± 0.01
BC	0.17 ± 0.01	0.07 ± 0.01	0.04 ± 0.00	0.03 ± 0.00

## 180 4 Discussion

181 We have demonstrated that HA outperforms various offline RL methods across three simulated object  
 182 rearrangement problems. We hope HA will inspire future work on how abstractions of independent  
 183 and symmetric entities can be learned and re-composed.

## References

- 184
- 185 [1] David Abel, D. Ellis Hershkowitz, Gabriel Barth-Maron, Stephen Brawner, Kevin O’Farrell,  
186 James MacGlashan, and Stefanie Tellex. Goal-based action priors. In *ICAPS*, pages 306–314.  
187 AAAI Press, 2015.
- 188 [2] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius  
189 Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan  
190 Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint*  
191 *arXiv:1806.01261*, 2018.
- 192 [3] Craig Boutilier, Richard Dearden, and Moises Goldszmidt. Exploiting structure in policy  
193 construction. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence*  
194 *- Volume 2, IJCAI’95*, pages 1104–1111, San Francisco, CA, USA, 1995. Morgan Kaufmann  
195 Publishers Inc.
- 196 [4] Craig Boutilier, Richard Dearden, and Moisés Goldszmidt. Stochastic dynamic programming  
197 with factored representations. *Artif. Intell.*, 121(1-2):49–107, 2000.
- 198 [5] Berk Calli, Arjun Singh, Aaron Walsman, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M  
199 Dollar. The ycb object and model set: Towards common benchmarks for manipulation research.  
200 In *2015 international conference on advanced robotics (ICAR)*, pages 510–517. IEEE, 2015.
- 201 [6] Michael Chang, Thomas L Griffiths, and Sergey Levine. Object representations as fixed  
202 points: Training iterative refinement algorithms with implicit differentiation. *arXiv preprint*  
203 *arXiv:2207.00787*, 2022.
- 204 [7] Michael B Chang, Tomer Ullman, Antonio Torralba, and Joshua B Tenenbaum. A compositional  
205 object-based approach to learning physical dynamics. *arXiv preprint arXiv:1612.00341*, 2016.
- 206 [8] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games,  
207 robotics and machine learning. 2016.
- 208 [9] Carlos Diuk, Andre Cohen, and Michael L. Littman. An object-oriented representation for  
209 efficient reinforcement learning. In *ICML*, volume 307 of *ACM International Conference*  
210 *Proceeding Series*, pages 240–247. ACM, 2008.
- 211 [10] Ben Eysenbach, Russ R Salakhutdinov, and Sergey Levine. Search on the Replay Buffer:  
212 Bridging Planning and Reinforcement Learning. In H. Wallach, H. Larochelle, A. Beygelzimer,  
213 F. d’ Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing*  
214 *Systems*, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.  
215 cc/paper/2019/file/5c48ff18e0a47baaf81d8b8ea51eec92-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/5c48ff18e0a47baaf81d8b8ea51eec92-Paper.pdf).
- 216 [11] Richard E. Fikes and Nils J. Nilsson. Strips: A new approach to the application of theorem  
217 proving to problem solving. *Artificial Intelligence*, 2(3):189–208, 1971. ISSN 0004-3702. doi:  
218 [https://doi.org/10.1016/0004-3702\(71\)90010-5](https://doi.org/10.1016/0004-3702(71)90010-5). URL [https://www.sciencedirect.com/  
219 science/article/pii/0004370271900105](https://www.sciencedirect.com/science/article/pii/0004370271900105).
- 220 [12] Natalia Gardiol and Leslie Kaelbling. Envelope-based Planning in Relational MDPs. In  
221 S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing*  
222 *Systems*, volume 16. MIT Press, 2003. URL [https://proceedings.neurips.cc/paper/  
223 2003/file/4a06d868d044c50af0cf9bc82d2fc19f-Paper.pdf](https://proceedings.neurips.cc/paper/2003/file/4a06d868d044c50af0cf9bc82d2fc19f-Paper.pdf).
- 224 [13] M. Ghallab, A. Howe, C. Knoblock, D. Mcdermott, A. Ram, M. Veloso, D. Weld, and D. Wilkins.  
225 PDDL—The Planning Domain Definition Language, 1998. URL [http://citeseerx.ist.  
226 psu.edu/viewdoc/summary?doi=10.1.1.37.212](http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.37.212).
- 227 [14] Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. Neural expectation maximization.  
228 *arXiv preprint arXiv:1708.03498*, 2017.
- 229 [15] Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Christopher Burgess,  
230 Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object repre-  
231 sentation learning with iterative variational inference. In *International Conference on Machine*  
232 *Learning*, pages 2424–2433. PMLR, 2019.

- 233 [16] Klaus Greff, Sjoerd van Steenkiste, and Jürgen Schmidhuber. On the binding problem in  
234 artificial neural networks. *arXiv preprint arXiv:2012.05208*, 2020.
- 235 [17] C. Guestrin, D. Koller, R. Parr, and S. Venkataraman. Efficient Solution Algorithms for Factored  
236 MDPs. *Journal of Artificial Intelligence Research*, 19:399–468, October 2003. ISSN 1076-  
237 9757. doi: 10.1613/jair.1000. URL [https://jair.org/index.php/jair/article/view/  
10341](https://jair.org/index.php/jair/article/view/10341). tex.ids: guestrin2003EfficientSolutionAlgorithms.
- 239 [18] Carlos Guestrin, Daphne Koller, Chris Gearhart, and Neal Kanodia. Generalizing plans to new  
240 environments in relational mdps. In *Proceedings of the 18th International Joint Conference  
241 on Artificial Intelligence, IJCAI’03*, page 1003–1010, San Francisco, CA, USA, 2003. Morgan  
242 Kaufmann Publishers Inc.
- 243 [19] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control:  
244 Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 245 [20] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model:  
246 Model-based policy optimization. *Advances in Neural Information Processing Systems*, 32,  
247 2019.
- 248 [21] Thomas Kipf, Gamaleldin F Elsayed, Aravindh Mahendran, Austin Stone, Sara Sabour, Georg  
249 Heigold, Rico Jonschkowski, Alexey Dosovitskiy, and Klaus Greff. Conditional object-centric  
250 learning from video. *arXiv preprint arXiv:2111.12594*, 2021.
- 251 [22] Ilya Kostrikov. JAXRL: Implementations of Reinforcement Learning algorithms in JAX, 10  
252 2021. URL <https://github.com/ikostrikov/jaxrl>.
- 253 [23] Tejas D Kulkarni, Ankush Gupta, Catalin Ionescu, Sebastian Borgeaud, Malcolm Reynolds,  
254 Andrew Zisserman, and Volodymyr Mnih. Unsupervised learning of object keypoints for  
255 perception and control. *NeurIPS*, 2019.
- 256 [24] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg  
257 Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with  
258 slot attention. *arXiv preprint arXiv:2006.15055*, 2020.
- 259 [25] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg  
260 Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with  
261 slot attention. In *NeurIPS*, 2020. URL [https://proceedings.neurips.cc/paper/2020/  
hash/8511df98c02ab60aea1b2356c013bc0f-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/8511df98c02ab60aea1b2356c013bc0f-Abstract.html).
- 263 [26] OpenAI. Robogym. <https://github.com/openai/robogym>, 2020.
- 264 [27] Luis Pineda, Brandon Amos, Amy Zhang, Nathan O. Lambert, and Roberto Calandra. Mbrl-  
265 lib: A modular library for model-based reinforcement learning. *Arxiv*, 2021. URL <https://arxiv.org/abs/2104.10159>.
- 267 [28] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language  
268 understanding by generative pre-training. 2018.
- 269 [29] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark  
270 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*,  
271 2021.
- 272 [30] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Si-  
273 mon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering  
274 atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 275 [31] Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate DALL-E learns to compose. *arXiv preprint  
276 arXiv:2110.11405*, 2021.
- 277 [32] Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate DALL-e learns to compose. In *Inter-  
278 national Conference on Learning Representations*, 2022. URL [https://openreview.net/  
279 forum?id=h00YV0We3oh](https://openreview.net/forum?id=h00YV0We3oh).

- 280 [33] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based  
281 control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages  
282 5026–5033. IEEE, 2012.
- 283 [34] Sjoerd Van Steenkiste, Michael Chang, Klaus Greff, and Jürgen Schmidhuber. Relational neural  
284 expectation maximization: Unsupervised discovery of objects and their interactions. *arXiv*  
285 *preprint arXiv:1802.10353*, 2018.
- 286 [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
287 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*  
288 *processing systems*, 30, 2017.
- 289 [36] Rishi Veerapaneni, John D Co-Reyes, Michael Chang, Michael Janner, Chelsea Finn, Jiajun Wu,  
290 Joshua Tenenbaum, and Sergey Levine. Entity abstraction in visual model-based reinforcement  
291 learning. In *Conference on Robot Learning*, pages 1439–1456. PMLR, 2020.
- 292 [37] C. Wang, S. Joshi, and R. Khardon. First order decision diagrams for relational mdps. *Journal*  
293 *of Artificial Intelligence Research*, 31:431–472, Mar 2008. ISSN 1076-9757. doi: 10.1613/jair.  
294 2489. URL <http://dx.doi.org/10.1613/jair.2489>.
- 295 [38] Ge Yang, Amy Zhang, Ari S. Morcos, Joelle Pineau, Pieter Abbeel, and Roberto Calandra.  
296 Plan2vec: Unsupervised Representation Learning by Latent Plans. In *Proceedings of The 2nd*  
297 *Annual Conference on Learning for Dynamics and Control*, volume 120 of *Proceedings of*  
298 *Machine Learning Research*, pages 1–12, 2020.
- 299 [39] Amy Zhang, Adam Lerer, Sainbayar Sukhbaatar, Rob Fergus, and Arthur Szlam. Composable  
300 Planning with Attributes. In *Proceedings of the 35th International Conference on Machine*  
301 *Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80,  
302 pages 5837–5846. JMLR.org, 2018.
- 303 [40] Daniel Zoran, Rishabh Kabra, Alexander Lerchner, and Danilo J Rezende. Parts: Unsupervised  
304 segmentation with slots, attention and independence maximization. In *Proceedings of the*  
305 *IEEE/CVF International Conference on Computer Vision*, pages 10439–10447, 2021.



## 306 A Additional Related Work

307 **Structured MDPs.** Many works have postulated that designing algorithms for structured MDPs will  
308 lead to improvements in sample efficiency and generalization over existing algorithms for the standard  
309 MDP formulation. The two closest types of MDP families to our work are factored MDPs [3, 4, 17],  
310 relational MDPs [37], and object-oriented MDPs [9]. Factored MDPs [3, 4, 17] assume that the  
311 environment can be represented by discrete attributes, and that transitions between these attributes  
312 can be modeled as a Bayesian network. Our work differs from these in that we do not assume access  
313 to these attributes and the dependency graph is also not assumed to be known. More importantly, the  
314 focus in this work is on using attributes to factorize a task into subcomponents; and in particular being  
315 able to generalize to new, more complex tasks at test time. Our approach is also related to Relational  
316 MDPs and Object-Oriented MDPs [1, 9, 12], where states are described as a set of objects, each of  
317 which is an instantiation of canonical classes, and each instantiated object has a set of attributes. Our  
318 work is especially related to [18], where the aim is to show that by using a relational representation of  
319 an MDP, a policy from one domain can generalize to a new domain through planning. However, we  
320 discover both objects (identities) and attributes (states), and achieve generalization through factorized  
321 planning, which grows linearly (rather than polynomially) in the number of factors.

## 322 B Implementation Details

323 This section details the implementation design decisions for each component of HA.

### 324 B.1 Level 1: abstracting visual features into sets of entities

325 Dynamic SLATE (dSLATE) maps a video demonstration to a sequence of transitions over sets of  
326 entities. It consists of two main components: SLATE [32] and a transformer [28, 35] dynamics model.  
327 Because SLATE itself also uses a transformer for generating image reconstructions, dSLATE uses  
328 two transformers: one to recover the observation model  $E$  and one to recover the dynamics model  $P$   
329 of the structured MDP. We use SLATE  $\prod_k \mathcal{H}^k \times \mathcal{O} \rightarrow \prod_k \mathcal{H}^k$  to infer entities  $\mathbf{h}_t$  from observation  
330  $o_t$  and initial guess  $\hat{\mathbf{h}}_t$ . The entity  $h^k$  is also referred to as a *slot* in [25, 32] and is split in half as  
331  $h^k = (z^k, s^k)$ . The first guess for each entity  $\hat{h}_1^k$  is sampled independently and identically distributed  
332 from a unit Gaussian, whose parameters are also trained. The hyperparameters are given in Tab. 5.

333 SLATE preprocesses the image with a discrete variational autoencoder [29] into a grid of image  
334 features, encodes these features into a grid of tokens, infers slots from this token grid with Slot  
335 Attention [25], which also produces an attention mask `attn` over the features each slot attends to.  
336 These slots are trained using a transformer decoder [28, 35] to autoregressively reconstruct the tokens  
337 using the slots as keys/values.

338 We use the dynamics model  $\prod_k \mathcal{S}^k \times \mathcal{A} \rightarrow \prod_k \mathcal{S}^k$  to predict a guess for  $\hat{\mathbf{s}}_{t+1}$  given the action  $a_t$  and  
339 the inferred entities  $\mathbf{s}_t$  from the previous time-step. This dynamics model is implemented also as a  
340 transformer decoder, taking the entity states as queries, and the entity states and action as keys/values.  
341 This enables us the dynamics model to model the future state of the slots as an equivariant function  
342 of how the action affects it and of how it interacts with other entity states.

343 We trained dSLATE on an offline dataset of 5000 video demonstrations of length five, with each  
344 frame transition showing one of four objects being moved to a different location. We used five  
345 slots, one more than the number of objects, following the convention used in Van Steenkiste et al.  
346 [34], Veerapaneni et al. [36].

### 347 B.2 Level 2: abstracting transitions over sets of entities into individual state transitions

348 HA constructs a graph of state transitions from a buffer of transitions over entity sets produced by  
349 dSLATE, as described in Alg. 1. Each transition records the identity  $z$ , state  $s$ , and `attn` of each  
350 entity of the entity sets  $\mathbf{h}_t$  and  $\mathbf{h}_{t+1}$ . We also record the *pre-condition* of the transition, which we  
351 explain further in Appdx. B.3. Both  $s$  or `attn` can serve as the representation of the state. We found  
352 that we obtained better clusterings when we used `attn` as the state for the *block-\** tasks and  $s$  as  
353 the state for the *robogym-rearrange* task. We also empirically found that certain choices of distance  
354 metric used for K-means clustering and binding (implemented as nearest-neighbors) depended

Number of epochs		200
Episodes per epoch		5K
Episode length		5
Batch size		32
Peak LR		(see caption)
LR warmup steps		30000
Dropout		0.1
Discrete VAE	Vocabulary Size	4096
	Temp. Cooldown	1.0 to 0.1
	Temp. Cooldown Steps	30000
	LR (no warmup)	0.0003
	Image Size	(see caption)
	Image Tokens	Image Size / 4
transformer decoder	Layers	4
	Heads	4
	Hidden Dim.	192
Slot attention	Slots	5
	Iterations	3
	Slot Heads	1
	Slot Dim. ( $h^k$ )	192
	Identity Dim. ( $z^k$ )	96
	State Dim. ( $s^k$ )	96
transformer dynamics	Layers	4
	Heads	4
	Hidden Dim.	96

Table 5: **Hyperparameters for training dSLATE** These hyperparameters are almost identical to those found in Singh et al. [32, Fig. 7], but because dSLATE operates on video demonstrations rather than static images, we changed some hyperparameters to save memory cost. We changed the batch size from 50 to 32, the number of transformer layers and heads from 8 to 4, the number of slot attention iterations from 7 to 3 without observing a significant change in performance. We used a peak learning rate of 0.0002 and an image size of 64 for *\*-rearrange*. We used a peak learning rate of 0.0003 and an image size of 96 for *block-stacking*.

355 on which choice of state representation we used, and this is summarized in Table 6. The K-means  
356 implementation is adapted from [https://github.com/overshiki/kmeans\\_pytorch](https://github.com/overshiki/kmeans_pytorch). We found  
357 that increasing the number of slot attention iterations improved the entities representations especially  
358 when generalizing to more numbers of objects, so even though we dSLATE trained with slot attention  
359 three iterations, for inferring the slots from the buffer we used seven iterations. Lastly, we found  
360 that the number of clusters used to for K-Means is the most important hyperparameter for creating  
361 a graph that reflected the state transitions. We swept over 16 to 50 clusters and report the optimal  
362 number of clusters we found in Table 7.

State representation	attn	$s$
<code>isolate</code> distance metric $d(\cdot, \cdot)$	cosine	cosine
<code>cluster</code> distance metric	IoU	squared Euclidean
<code>bind</code> distance metric	cosine	squared Euclidean

Table 6: **Hyperparameters for constructing the transition graph with HA**

	<i>block-rearrange</i>	<i>robogym-rearrange</i>	<i>block-stacking</i>
number of clusters	30	45	47

Table 7: **Number of clusters used for constructing the nodes of the transition graph.**

363 **B.3 Action selection**

364 Using dSLATE and the transition graph from HA, HA returns which action to execute in the  
 365 environment given a goal and current observation, as described in Alg. 2. It first infers goal constraints  
 366  $\mathbf{h}_g$  and current entities  $\mathbf{h}_t$  from the goal observation  $o_g$  and current observation  $o_t$ . It then uses the  
 367 `align` procedure to align the indices of the entities in  $\mathbf{h}_g$  and  $\mathbf{h}_t$  and uses the `select-constraint`  
 368 to choose the index  $k$  of the entity to affect. It binds  $h_t^k$  and  $h_g^k$  to the graph and returns the action  
 369 associated with the edge between their respective nodes. Because HA is a non-parameteric method, it  
 370 could be the case the graph does not contain such an edge. In this case, we sample a random action in  
 371 the environment, but future work will replace this step with a more sophisticated method.

372 To implement `align` we use the `scipy.optimize.linear_sum_assignment` implementation of  
 373 the Hungarian algorithm, with Euclidean distances between the  $z^k$ 's as the matching cost.

374 In `select-constraint`, we are given the set of current entities  $\mathbf{h}_t$  whose indices are aligned with the  
 375 goal constraints  $\mathbf{h}_g$  and returns the index  $k$  of the goal constraint to satisfy next. By HA' construction,  
 376 the edge between the nodes that  $h_t^k$  and  $h_g^k$  are bound to is the state transition that would be executed  
 377 if the action associated to the edge were taken in the environment. The `select-constraint`  
 378 procedure consists of three steps: (1) filtering possible transitions from impossible transitions (2)  
 379 ranking transitions (3) sampling a transition.

380 **Filtering** The filtering step implements HA' model of possibility and impossibility. In the filtering  
 381 step, we consider, for each  $k$ , the transition between the nodes that  $h_t^k$  and  $h_g^k$  are bound to and mark  
 382 the transition as possible or impossible. It then returns the indices  $k$  over  $\mathbf{h}_t$  and  $\mathbf{h}_g$  whose associated  
 383 transition from  $h_t^k$  and  $h_g^k$  is possible.

384 According to HA, a state transition between node  $[i]$  and node  $[j]$  is *possible* if its preconditions  
 385 are met and there exists an edge for that state transition in the graph, and *impossible* otherwise. An  
 386 intuitive example of an impossible transition is to place a block midair at some intended height, but  
 387 this transition becomes possible if prior to the transition there already exists a stack of blocks that  
 388 would support the block if the block were to be placed at that intended height. The existence of  
 389 this supporting stack is thus the *precondition* for the transition to occur, rendering the stacked block  
 390 *dependent* on the blocks supporting it.

391 When there are no dependencies among the entities, as in the *\*-rearrange* tasks where any object can  
 392 be moved to any open location without considering where other objects are, any transition present in  
 393 the graph is possible. When there are dependencies among the entities, as in *block-stacking*, we take  
 394 the precondition of the transition into account. Although a precondition of a transition from node  $[i]$   
 395 to node  $[j]$  could be a function of both the source node  $[i]$  and destination node  $[j]$ , for simplicity in  
 396 this paper we consider preconditions as only a function of the destination node  $[j]$ , which rules out  
 397 the possibility of placing a block in midair like the above example.

398 Because a precondition in general is a set of constraints that need to be satisfied for the transition to  
 399 be possible, we represent the precondition of a transition into node  $[j]$  (denoting the state  $s_*^{[j]}$ ) as the  
 400 set of context states  $s_*^{[j']}$  that are always present whenever the state  $s_*^{[j]}$  is present. Concretely, a block  
 401 at height 3 at location  $x$  is always accompanied by the presence of some block at height 2 and some  
 402 block at height 1, both at location  $x$ ; the states denoting (location  $x$ , height 2) and (location  $x$ , height  
 403 1) are the context states of the state (location  $x$ , height 3) and therefore serves as its precondition. We  
 404 thus implement the precondition of a transition into node  $[j]$  by recording the indices  $j'$  of the nodes  
 405 of states that are always present when node  $[j]$  is a destination node. To test whether a precondition is  
 406 satisfied for a given scene, we check if all nodes in the precondition set have a corresponding concrete  
 407 entity that can be bound to it.

408 **Ranking** The filtering step removes the indices from the entities  $\mathbf{h}_t$  and goal constraints  $\mathbf{h}_g$  whose  
 409 transitions are impossible, yielding a possibly smaller set of entities  $\tilde{\mathbf{h}}_t$  and constraints  $\tilde{\mathbf{h}}_g$ . That is, if  
 410  $|\mathbf{h}_t| = |\mathbf{h}_g| = K$ , then  $|\tilde{\mathbf{h}}_t| = |\tilde{\mathbf{h}}_g| = \tilde{K} \leq K$ .

411 The goal of the ranking step is to compute a ranking among the indices of  $\tilde{\mathbf{h}}_t$  and  $\tilde{\mathbf{h}}_g$  for choosing  
 412 which index  $k$  to actually select to affect with an action. Intuitively, we should rank indices  $k$

413 according to how different  $s_t^k$  and  $s_g^k$  are because a difference indicates that the constraint  $h_g^k$  is not  
 414 satisfied. We reuse the distance metric  $d(\cdot, \cdot)$  used for `isolate` to implement this ranking.

415 **Sampling** The goal of the sampling step is to select a  $k \in \{1, \dots, \tilde{K}\}$  whose associated entity  
 416 we will affect with an action, given our ranking. One way to do this is to simply choose  $k$  as  
 417  $k = \arg \max_{k' \in \{1, \dots, \tilde{K}\}} d(s_t^{k'}, s_{t+1}^{k'})$  as in `isolate`, but we empirically found that sampling  $k$  from  
 418 a Categorical distribution whose pre-normalized probabilities are given by  $d(s_t^{k'}, s_{t+1}^{k'})$  resulted in  
 419 better task performance so we used this stochastic sampling approach. One explanation for why using  
 420 the argmax may be worse is that it relies on the distance metric  $d(\cdot, \cdot)$ , and the state representation  $s$ ,  
 421 to be such that the distance metric flawlessly assigns high value to entities  $k$  that need to be moved  
 422 and low value to entities  $k$  that do not need to be moved. But because the state space  $\mathcal{S}$  is learned  
 423 through the dSLATE training process without explicit supervision on the geometry of the space, a  
 424 pair of points that should be farther apart than another set of points may not be accurately reflected  
 425 by using a fixed distance metric  $d(\cdot, \cdot)$ . Future work will investigate imposing explicit supervision on  
 426 the geometry of  $\mathcal{S}$ .

## 427 C Baseline Implementation Details

428 **Random (Rand)** The random policy takes actions using `env.action_space.sample()`.

429 **Behavior cloning (BC)** This approach trains a policy to output the actions directly taken in the  
 430 provided dataset. We use an MSE loss to train the policy to imitate the actions.

431 **Implicit Q-learning (IQL)** IQL is a simple, offline RL approach that uses temporal difference  
 432 (TD) learning with the dataset actions and trains a behavior policy value function. To produce an  
 433 optimal value function, IQL estimates the maximum of the Q-function using expectile regression  
 434 with an asymmetric MSE using the following objectives:

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^\tau(Q_{\hat{\theta}}(s, a) - V_\psi(s))] \text{ where } L_2^\tau(u) = |\tau - \mathbb{1}(u < 0)|u^2 \quad (1)$$

$$L_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} [(r(s, a) + \gamma V_\psi(s') - Q_\theta(s, a))^2] \quad (2)$$

$$L_\pi(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [\exp(\beta(Q_{\hat{\theta}}(s, a) - V_\psi(s))) \log \pi_\phi(a|s)]. \quad (3)$$

435 The  $V(s)$  estimates are used for TD-backups and the optimal policy is extracted with advantage-  
 436 weighted behavioral cloning.

437 **Model predictive control (MPC)** This approach uses model predictive control with the cross  
 438 entropy method (CEM) to select actions, using the transformer dynamics model of dSLATE to  
 439 perform rollouts in latent space. This is similar to the approached used in OP3 [36], except that we  
 440 use more recently proposed architectural components (slot attention [25] instead of IODINE [15], a  
 441 transformer instead of a graph network [2, 7, 34]) so our MPC results are not directly comparable to  
 442 that of OP3. We use the same dSLATE checkpoint that was used for HA.

443 We implement this MPC baseline using the `mbr1-lib` library [27] with 10 CEM iterations, an elite  
 444 ratio of 0.05, and a population size of 250 which was the best configuration we found that fit within a  
 445 wall clock budget of two days for 8 objects and 100 test episodes. We swept over CEM iterations of  
 446  $[5, 10, 20]$ , elite ratio of  $[0.05, 0.1, 0.2]$ , and population sizes of  $[250, 500, 1000]$ , and found that the  
 447 elite ratio was the most important hyperparameter.

448 The cost function is computed by first aligning the predicted slots  $\mathbf{h}_T$  and goal constraints  $\mathbf{h}_g$  using  
 449 the same `align` procedure in B.3, and then adding up the squared Euclidean distance between slots  
 450 as  $cost = \sum_k (h_T^k - h_g^k)^2$ .

451 **Monolithic graph search (MGS)** This approach is an ablation to HA that does not construct a  
 452 graph over state transitions of individual entities but instead constructs a graph over state transition  
 453 over entity sets, i.e. each transition is  $(s, a, s')$  rather than  $(s^k, a, s^{k'})$ . As with MPC, we use the  
 454 same dSLATE checkpoint that was used for HA.

455 The purpose of this ablation is to elucidate the benefit of factorizing the transition graph over entities  
 456 rather than entity sets. Because nodes in this transition graph represent a set of entity states rather



Figure 5: An example of solving a task in the robogym rearrange environment used in this paper.

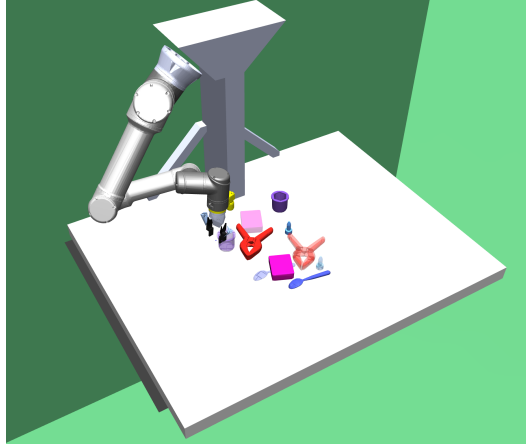


Figure 6: The original Robogym rearrange setup

457 than individual entity states, we use Dijkstra’s algorithm, as in [10, 38, 39] to plan a unbroken path  
 458 from the node the initial observation is bound to to the node a goal observation is bound to. For each  
 459 time-step, we plan a path along the nodes using Dijkstra’s algorithm, then return the action associated  
 460 with the first edge along that path. Like HA, MGS is a non-parametric model, which means that for a  
 461 set of entities to be bound to a node in the graph, that node must contain the exact set of entity states  
 462 corresponding to the states of the entities. If we do not successfully bind to the graph, or if we do not  
 463 find a path between the current node and the goal node, we sample a random action as HA does.

## 464 D Environment Details

465 *Block-rearrange* is our simplest environment, and *robogym-rearrange* and *block-stack* each add  
 466 complexity along different axes. States and identities correspond to object locations and appearance  
 467 respectively. In *block-rearrange* (Fig. 2a), all objects are the same size, shape, and orientation.  
 468  $\mathcal{S}$  covers 16 locations in a grid.  $\mathcal{Z}$  is the continuous space of red-green-blue values from 0 to 1.  
 469 *robogym-rearrange* (fig. 2b) adapts the rearrange environment from OpenAI [26] and removes the  
 470 assumption from *block-rearrange* that all objects are the same size, shape, orientation, and the  
 471 assumption of predefined locations. *block-stack* (fig. 2c) adds preconditions and postconditions on  
 472 whether objects can be moved: blocks can only be picked if from the top of a stack, and blocks  
 473 can only be placed at a given height if there is an object beneath to support it (otherwise it falls).  
 474 *Block-rearrange* and *Block-stack* are implemented in PyBullet [8] while *Robogym-rearrange* is  
 475 implemented in Mujoco [33].

476 **Environments** In *block-rearrange*, all objects are the same size, shape, and orientation and can  
 477 exist in any one of 16 locations in a grid. Colors are sampled in a continuous space of red-green-blue  
 478 values in  $[0, 1]$ .

479 The *robogym-rearrange* environment (see figures 5 and 6) is adapted from the rearrange environment  
 480 in OpenAI’s Robogym simulation framework [26] and removes the assumption from *block-rearrange*  
 481 that all objects are the same size, shape, and orientation and the assumption of predefined locations.  
 482 Furthermore, due to 3D perspective, the objects can look slightly different in different locations. The  
 483 objects are uniformly sampled from a set of 94 meshes consisting of the YCB object set [5] and a  
 484 set of basic geometric shapes, with colors sampled from a set of 13. The camera angle is a bird’s  
 485 eye view over the table, and the size of each object is normalized by its longest dimension, so tall

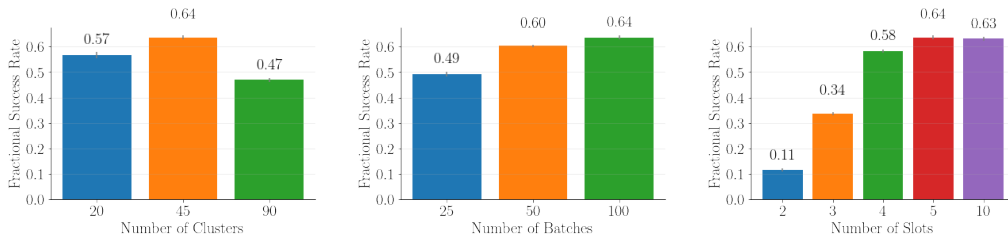


Figure 7: The performance of our method as the number of initialized clusters and batches from the training set used to construct the graph, and the number of slots are varied.

486 thin objects appear smaller. The objects’ target positions are randomly sampled such that they don’t  
 487 overlap with each other or any of the initial positions, and the target orientation is set to be unchanged.  
 488 Due to the continuous nature of this environment, we define a match threshold of at most 0.05 for  
 489 both the initial pick position and the goal placement (the table dimensions are 0.6 by 0.8).

490 The *block-stack* environment adds preconditions and postconditions on whether objects can be moved:  
 491 blocks can only be picked if they are at the top of a stack, and blocks can only be placed at a given  
 492 height if there is already an object beneath to support it (otherwise it falls).

493 **Sensorimotor interface** Each observation is a tuple of an initial image displaying the current  
 494 observation and a goal image displaying constraints to be satisfied – the goal locations of the objects.  
 495 Each action is a tuple  $(w, \Delta w)$ , where  $w$  is a three-dimensional Cartesian coordinate  $(x, y, z)$  in the  
 496 environment arena.

497 For the *\*-rearrange* tasks, objects are initialized at random non-overlapping locations that also do not  
 498 overlap with their goal locations. For these tasks the  $z$  (height) coordinate is always fixed. For the  
 499 *block-stack* task, the goal locations are generated by randomly picking objects from the tops of stacks  
 500 and placing them on other stacks. For this task the  $y$  (depth) coordinate is always fixed.

501 An object is picked if  $w$  is within a certain threshold of its location. For *block-\** where object locations  
 502 are fixed points in a grid the object is snapped to the nearest grid location to  $w + \Delta w$ . Constraints  
 503 are considered satisfied if objects are placed within a certain threshold of their target location.

## 504 E Additional Results

505 This section presents additional results and analyses of HA. We first analyze the sensitivity of task  
 506 performance to several hyperparameters used in HA from creating the graph: the number of clusters,  
 507 number of buffer size, and the number of slots used in slot attention. We next tested whether giving  
 508 our model-based baselines more computation time compared to HA would improve their performance  
 509 to be comparable to HA’s.

### 510 E.1 Analysis of key hyperparameters

511 In this section, we analyze the sensitivity of our method to various hyperparameters, evaluated in  
 512 the robogym environment with four objects in the complete goal specification. As Fig. 7 shows,  
 513 performance depends on the number of initialized clusters and the number of batches from the  
 514 training set used to construct the graph. With too few clusters, the clusters are too coarse-grained  
 515 to differentiate objects in significantly different positions, while with too many the performance  
 516 deteriorates as the data is needlessly split into duplicate clusters. Performance improves with more  
 517 data, as the graph has better coverage. Although our model performs worse when there are insufficient  
 518 slots to represent all objects present in the environment plus one empty slot, performance is barely  
 519 impacted by having double the number of necessary slots. Our method can thus still work in  
 520 environments with an unknown but upper-bounded number of objects.

### 521 E.2 More challenging evaluation settings

522 We analyzed HA in more challenging settings that reflect the noisy nature of real-world robotics. As  
 523 Fig. 8 (left) shows, HA is robust to the addition of Gaussian noise to the action at every time step,  
 524 up until the noise variance is comparable to the maximum distance for successful picking and goal  
 525 placements. The performance also remains high given significantly fewer steps (Fig. 8, right).

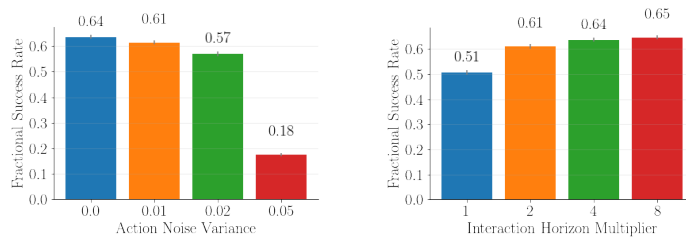


Figure 8: The performance of HA on *robogym-rearrange* as we vary the amount of noise added to the actions, and the interaction horizon (as a multiple of the minimum steps needed to complete the task).

### 526 E.3 More computation time for model-based baselines

527 We tested whether doubling the computation time for the model-based baselines would improve their  
 528 performance to be comparable to HA’s. For the results in the main paper, we capped the length of the  
 529 episode as 4x the minimum number of actions required to solve the task. In Fig. 9, we vary this  
 530 interaction horizon multiplier from 1x to 8x. HA degrades less with shorter interaction horizons  
 531 compared to the baselines. We find that MGS performs similar to the random baseline. Since MGS  
 532 takes a random action if it cannot bind the given entity set to its graph, this result suggests that the  
 533 space of subsets of entities is so combinatorially large that MGS does not successfully bind to the  
 534 graph most of the time. We verified that this is the case by inspecting when MGS takes random  
 535 actions. MPC performs the worst out of all the methods, performing worse than random. We tested  
 536 that the cost function described in C ranks latents that match the goal constraint with lower cost than  
 537 randomly sampled latents, which suggests that the main source of error is due to the inaccuracy in the  
 538 prediction rollouts. This can be expected, as learned models suffer from compounding errors when  
 539 rolled out [20] and prior methods that use MPC for object-centric methods only roll out for very short  
 540 horizons [36].

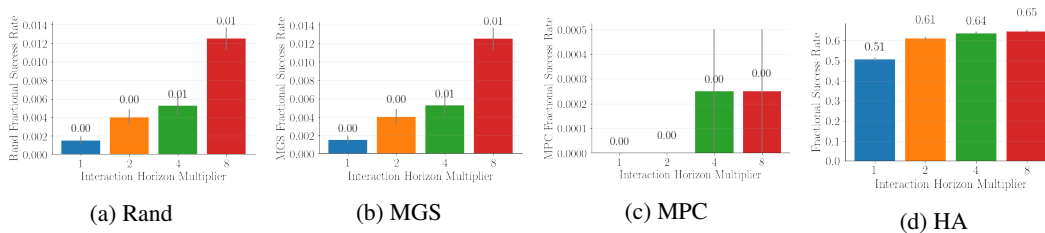


Figure 9: **Varying interaction horizon.** The performance of the MGS (b) and MPC (c) baselines compared to HA (d, reproduced from Fig. 8) and the random baseline (a) on *robogym-rearrange* as we vary the interaction horizon (as a multiple of the minimum steps needed to complete the task). *Note that the scale of the y-axis is not the same.* While a longer horizon improves performance, HA still achieves at least 50x better accuracy with an interaction horizon multiplier of 1 than the performance obtained by increasing the interaction horizon multiplier for the model-based baselines to 8.