

BINDING ACTIONS TO OBJECTS IN WORLD MODELS

Ondrej Biza^{1,*}, Robert Platt^{1,*}, Jan-Willem van de Meent^{1,2,*},
Lawson L. S. Wong^{1,*} and Thomas Kipf³

¹ Northeastern University, Boston, MA, USA

² University of Amsterdam, Netherlands

³ Google Research, Brain Team

* Equal contribution. Correspondence to: biza.o@northeastern.edu.

ABSTRACT

We study the problem of binding actions to objects in object-factored world models using action-attention mechanisms. We propose two attention mechanisms for binding actions to objects, soft attention and hard attention, which we evaluate in the context of structured world models for five environments. Our experiments show that hard attention helps contrastively-trained structured world models to learn to separate individual objects in an object-based grid-world environment. Further, we show that soft attention increases performance of factored world models trained on a robotic manipulation task. The learned action attention weights can be used to interpret the factored world model as the attention focuses on the manipulated object in the environment. Source code: https://github.com/ondrejba/action_attention_iclr_22.

1 INTRODUCTION

Reinforcement learning agents often operate in domains that contain multiple objects, such as robotic manipulation (e.g. Janner et al. (2019); Li et al. (2020)) or game playing (e.g. Bellemare et al. (2013); Guss et al. (2019)). An action performed by an agent in an object-based environment is unlikely to affect all objects at once. A household robot might pick up a single object or push several objects, but it will not move all objects in a room at once. In Atari Pong, the agent only controls the paddle while the other two objects—the ball and the other player—are not directly influenced by actions. The sparsity of interactions between actions and objects is a useful inductive bias for an agent’s world model.

We implement a *learned* binding of an action to objects, which is constrained to be sparse. The sparsity is enforced by an attention mechanism over objects, with individual weights being constrained to sum to one using the softmax operator. Prior work considered structured action spaces in world models for grid-worlds with objects (Kipf et al., 2020), a Tetris-like block stacking task (Janner et al., 2019; Veerapaneni et al., 2019) and in model-free learning for 2D construction from blocks (Bapst et al., 2019). Unlike our work, the binding of an action to a particular object, or an edge between objects in a graph, was hand-coded, without any learning. Goyal et al. (2021) proposed an attention mechanism for binding production rules (Lovett & R. Anderson, 2005), which include information about action selection, to slots in a structured model.

We propose two different attention mechanisms and test them in object-based grid worlds, two Atari games and a simulated robotic manipulation task. Firstly, we work with contrastively-learned structured world models (C-SWMs, Kipf et al. (2020)) that jointly learn to factor states and predict the dynamics of the environment. This class of models often requires a strong bias in order to learn to properly distinguish objects without explicit supervision. To this end, we encode the bias that an action affects exactly one object using a categorical random variable. We call this approach *hard attention*, drawing on computer vision literature (Sharma et al., 2015). Secondly, we test our idea in a realistic robotic manipulation domain, where the individual states are already factored into object and our task is to learn factored world models (FWMs, Biza et al. (2022)). Here, the robot is tasked with building towers from blocks; all blocks in a tower might potentially be affected by a robot’s action. Hence, the categorical action binding bias is too restrictive. Instead, we propose a *soft attention* module, which transforms and weights actions for each object based on their predicted

impact on said object. It is based on single-head self-attention (Vaswani et al., 2017). We limit our exploration of action binding to C-SWMs and FWMs; see Battaglia et al. (2018); Greff et al. (2020) for a review of other structured model architectures.

Our results indicate two positive outcomes and one negative outcome. When each object can be individually affected by an action, hard action attention proves an excellent inductive bias for learning to factor states. Conversely, when some objects are not directly affected by actions, neither soft nor hard action attention makes a difference. Finally, we show that soft attention can lead to large improvements in a pick-and-place robotic manipulation task.

In summary, our main contributions are:

#1: We propose two novel attention modules over actions, soft and hard attention, that improve the ability of structured models to factor states and to predict transition dynamics.

#2: We identify the phenomenon wherein a structured world model captures highly correlated information in all its available slots. This happens in Atari Pong and Space Invaders, where only the player-controlled paddle or spaceship is directly affected by actions.

2 PRELIMINARIES: STRUCTURED AND FACTORED WORLD MODELS

World models are trained to make predictions about the state of an environment s^0 after applying a sequence of actions a^1, a^2, \dots, a^T . It is commonly assumed that the state is represented as an image and the transition dynamics are deterministic. We make no assumptions about the action space. World models used in this paper learn a latent space without image reconstruction. That is, the model encodes s^0 into a latent encoding z^0 and predicts the future $z^{1:T}$ based on $a^{1:T}$.

Structured world models encode an image s into a latent state composed of K slots $z_{1:K}$, each represented as a feature vector. The model learns to assign objects in s to individual slots (i.e. to factor the scene) without any direct supervision. We use the contrastively-trained structured world model (C-SWM, Kipf et al. (2020))—it learns structured embeddings using a convolutional encoder E_ϕ and predicts the effects of actions using a graph neural network (GNN) T_θ (Gori et al., 2005). The GNN uses a f_{node} and an edge network f_{edge} to compute the next latent state of the k th object:

$$\hat{z}_k^{t+1} = z_k^t + T_\theta(z^t, a^t)_k = z_k^t + f_{\text{node}}\left(z_k^t, a^t, \sum_{i \neq k} f_{\text{edge}}(z_i^t, z_k^t)\right). \quad (1)$$

The model is trained by contrastive learning using the following loss:

$$L(z^t, \bar{z}, z^{t+1}, \hat{z}^{t+1}) = \sum_k \left\| z_k^{t+1} - \hat{z}_k^{t+1} \right\|_2^2 - \max \left\{ 0, \gamma - \sum_k \left\| z_k^t - \bar{z}_k \right\|_2^2 \right\}. \quad (2)$$

The negative example \bar{z} is an encoding of a state randomly sampled from the training set and γ is a hyper-parameter. Intuitively, the negative term in the loss states that a pair of randomly sampled states should be at least γ apart in the embedding space. This prevents the model from converging to a trivial solution.

We use **factored world models** to refer to world models that operate over factored states. That is, they learn to model the transition dynamics of individual factors without having to learn to discover factors. Biza et al. (2022) proposed a factored world model for pick-and-place robotic manipulation tasks. The model receives a factored state $s_{1:K}$ composed of images of the K object present in a scene. Each state factor is individually encoded into a latent factor: $z_{1:K} = \langle f_{\text{enc}}(s_1), f_{\text{enc}}(s_2), \dots, f_{\text{enc}}(s_K) \rangle$. The transition model and the training loss follow C-SWM with the exception of two design decisions in the transition model. (1) several layers of GNNs are used and skip connections are added in between layers and (2) the action is provided not only to the node network f_{node} but also to the edge network f_{edge} . We experiment with different numbers of GNN layers, but omit the edge actions.

3 BINDING ACTIONS TO OBJECTS

In this section, we describe two attention mechanisms that bind actions to object slots. In both cases, we start with a factored latent state $z_{1:K}$, which is an output of an encoder of either C-SWM or

FWM. We want to compare the individual factors to the action in order to determine which slots are affected. First, we transform individual factors into key vectors, $k = \langle j_k(z_1), j_k(z_2), \dots, j_k(z_K) \rangle$, and the action into a query vector $q = j_q(a)$. Both j_k and j_q are Multi-Layer Perceptrons. The keys and queries are used to compute attention weights α over object slots 1 to K for action a :

$$\alpha = \text{softmax}(k_1^T q, k_2^T q, \dots, k_K^T q). \quad (3)$$

It is important to note that in both the hard and soft attention modules described below, the attention weights end up transforming the action that is given to the node network of the GNN transition model. For the k th object slot, instead of computing $f_{\text{node}}(z_k, a, \dots)$ we will use $f_{\text{node}}(z_k, a'_k, \dots)$, where a'_k is an action constructed specifically for the k th object slot.

Hard attention: We define a categorical variable M that represents the membership of an action to a particular factor:

$$M \sim \text{Categorical}(\alpha_1, \alpha_2, \dots, \alpha_K). \quad (4)$$

This random variable only plays a role when predicting the transitions and computing the positive term of the contrastive loss (the negative term only depends on the current encoded state). During training, we compute the expectation of M via summation over all of its possible assignments (same as the number of object slots / factors). During inference, we take the assignment with the highest probability. The random variable is used as follows:

$$\hat{z}_k^{t+1} = z_k^t + T_\theta(z^t, \text{pad}(a^t, m))_k \mid M = m, \quad (5)$$

with the $\text{pad}(a^t, m)$ operator assigning action a^t to the m th slot and zeros to all other slots. Finally, the positive term of the contrastive loss function is computed as follows:

$$\mathbb{E}_M \left[\left\| z_k^{t+1} - \hat{z}_k^{t+1} \right\|_2^2 \right] = \sum_{m=1}^K \alpha_m \left\| z_k^{t+1} - [z_k^t + T_\theta(z^t, \text{pad}(a^t, m))_k] \right\|_2^2. \quad (6)$$

Soft attention: Unlike hard attention, soft attention weights individual actions provided to slots by attention weights without defining a probabilistic model. We follow the single-head self-attention module from the Transformer (Vaswani et al., 2017). The action a is first transformed into a value vector $v = j_v(a)$. Then, the value vector is multiplied by the previously computed attention weights for each factor: $a' = \langle \alpha_1 v, \alpha_2 v, \dots, \alpha_K v \rangle$. We then pass the (newly) factored action a' into the transition model. Note the individual factors in the action are all equal up to a multiplicative factor. Table 2 in the Appendix visualizes α in a robotic pick-and-place task.

4 EXPERIMENTS

We aim to answer the following questions:

1. Does hard attention help C-SWM factor environments where each object can be independently affected by an action (2D Shapes and 3D Cubes)?
2. Does the same apply when some objects are only indirectly (or not at all) affected by actions (Atari Pong and Space Invaders)?
3. Does soft attention increase the performance of factored world models in robotic manipulation?

4.1 ENVIRONMENTS

We use two grid-world (2D Shapes and 3D Cubes) and two Atari environments (Atari Pong and Space Invaders) from Kipf et al. (2020) and a robotic manipulation environment from Biza et al. (2022). We make changes to 2D Shapes and 3D Cubes in order to better study the effect of action attention. We describe each environment below.

- **2D Shapes and 3D Cubes:** In both environments, five objects of different shapes and colors are placed in a 5×5 grid-world. Each object can be moved independently and the objects are not allowed to pass through each other. Crucially, Kipf et al. (2020) represented actions as a pair of object index and direction of movement. As the action space was factored by design, C-SWM could use this information to make state factorization easier. We change the action space to represent object position and direction of movement. Hence, the model does not know the index of the object being moved, making state factorization more challenging.

Table 1: Results for 2D grid-world *Shapes* and 3D grid-world *Cubes* with unfactored states and actions (different from Kipf et al. (2020), where actions are factored) as well as Atari Pong and Space Invader datasets from Biza et al. (2021). We report means and standard deviations for Hits@1 and Mean Reciprocal Rank for 1, 5 and 10 step predictions. Each model was run 20 times in *Shapes* and *Cubes*, and 4 times in Atari. The results for Atari differ from Biza et al. (2021) because we use 1k instead of a 100 negative states during evaluation, making the task harder. Additionally, we report correlation between individual slots of C-SWM.

	Model	1 Step		5 Steps		10 Steps		Slot
		H@1	MRR	H@1	MRR	H@1	MRR	Correl.
2D GRID SHAPES	C-SWM	28.7 \pm 2.6	42.9 \pm 2.6	3.7 \pm 0.8	10.9 \pm 1.6	1.5 \pm 0.4	5.4 \pm 1.0	1.00 \pm 0.00
	Soft Attention	30.9 \pm 2.9	45.7 \pm 2.6	4.3 \pm 0.8	12.5 \pm 1.5	1.7 \pm 0.4	6.3 \pm 1.0	1.00 \pm 0.00
	Hard Attention	98.7 \pm 5.3	99.1 \pm 3.6	95.5 \pm 16.9	96.3 \pm 14.5	93.4 \pm 20.5	94.3 \pm 18.7	0.08 \pm 0.08
3D GRID CUBES	C-SWM	32.0 \pm 3.4	46.3 \pm 3.1	4.8 \pm 1.2	13.3 \pm 2.1	2.0 \pm 0.5	7.0 \pm 1.4	1.00 \pm 0.00
	Soft Attention	40.8 \pm 12.1	54.5 \pm 10.6	9.0 \pm 6.3	20.0 \pm 8.4	4.0 \pm 2.9	11.3 \pm 5.1	0.99 \pm 0.02
	Hard Attention	97.3 \pm 5.9	98.3 \pm 3.9	88.8 \pm 18.2	91.6 \pm 14.3	79.2 \pm 24.5	83.3 \pm 20.8	0.21 \pm 0.18
ATARI PONG	C-SWM	93.2 \pm 1.1	96.3 \pm 0.6	36.3 \pm 3.8	49.5 \pm 3.8	11.2 \pm 2.4	19.5 \pm 1.8	0.84 \pm 0.11
	Soft Attention	93.8 \pm 0.8	96.6 \pm 0.4	26.1 \pm 6.3	38.5 \pm 6.5	7.3 \pm 4.3	13.0 \pm 5.6	0.87 \pm 0.10
	Hard Attention	92.9 \pm 0.2	96.2 \pm 0.1	28.3 \pm 7.1	41.9 \pm 6.9	9.6 \pm 3.6	16.8 \pm 4.7	0.85 \pm 0.10
ATARI SPACE	C-SWM	97.1 \pm 0.2	98.5 \pm 0.1	86.0 \pm 0.9	91.7 \pm 0.5	74.9 \pm 1.7	83.4 \pm 1.0	0.82 \pm 0.13
	Soft Attention	96.7 \pm 0.6	98.3 \pm 0.3	86.2 \pm 0.8	91.7 \pm 0.5	76.2 \pm 0.7	84.0 \pm 0.3	0.67 \pm 0.05
	Hard Attention	96.8 \pm 0.7	98.3 \pm 0.3	84.0 \pm 1.3	90.2 \pm 0.8	73.9 \pm 2.0	82.3 \pm 1.2	0.80 \pm 0.13

- **Atari Pong and Space Invaders:** C-SWM is trained and tested on modeling short action sequences (10 timesteps) in these two games. We use the approach from Biza et al. (2021), where a trained agent acts for a random number of steps to reach an interesting starting state from which a random policy is rolled out. The collected datasets contain only the random rollouts so that the agent experiences both good and bad actions.
- **Robotic Pick-and-Place:** A simulated UR5 robotic arm manipulates six cubes. The action space is the continuous space of (x, y) coordinates, where a pre-programmed `pick` or `place` action is executed. The environment is observed by two side-viewing RGB cameras and additional two RGB cameras capture the content of the robot’s hand. These observations are post-processed to extract individual objects and their bounding boxes to capture each object independently (i.e. a factored state space). The training task is to make a tower of four cubes; the model is then expected to transfer to four testing tasks—wall, stairs, two towers of three and three towers of two—without additional training.

4.2 SHAPES, CUBES AND ATARI GAMES

Setup: The C-SWMs are evaluated on their ability to predict future latent states in an evaluation dataset using ranking metrics (Hits@1 and MRR, see Appendix A.1). These metrics check if the predicted latent state is closer to the encoding of the true next state compared to a set of negative states. Intuitively, we both test the ability of the model to predict the future and the ability of the encoder to distinguish between states.

Result: Ranking scores are reported in Table 1. We would like to highlight the performance of hard action attention in 2D *Shapes* and 3D *Cubes*. In these environments, baseline C-SWM fails because the action space is not factored per individual object. In contrast, C-SWM with hard attention encodes the correct inductive bias and learns to assign individual objects to separate slots in most cases.

We do not see the same increase in performance in Atari games: baseline C-SWM, soft attention and hard attention achieve comparable results. To explain this result, we introduce a new metric to measure the correlation between individual slots. If all slots capture the same information given a single input state, the correlation coefficient is going to be 1 (please see Appendix A.1). We see that successful models in 2D *Shapes* and 3D *Cubes* have a low slot correlation, as each individual object should be represented by a different slot. Conversely, all models in Atari (as well

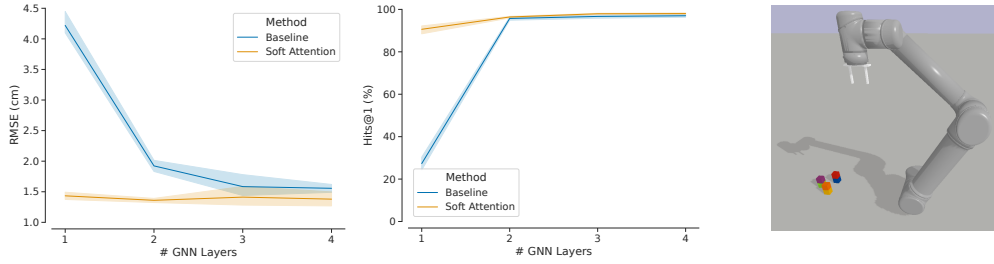


Figure 1: Zero-shot transfer results in block position prediction error (left, the lower the better) and action action sequence ranking (middle, the higher the better) in a UR5 pick-and-place environment with six cubes (right). We average over four random seeds and shade the 95% confidence interval.

as failed models in 2D Shapes and 3D Cubes) have high correlation between slots. We show one example of this phenomenon in Figure 2 in the Appendix where the latent space of all slots is nearly identical (the first slot seems to be flipped, but follows the same pattern).

The correlation metric results suggest that C-SWM has no need for more than one slot in Atari games. This is supported by the original results from Kipf et al. (2020). We hypothesize that since the action often affects only one object (the paddle in Pong and the spaceship in Space Invaders), there is not enough information to guide C-SWM to factor objects not directly affected by the action.

4.3 ROBOTIC PICK-AND-PLACE

Setup: We replicate the experiment from Biza et al. (2022), where a UR5 robotic arm picks and places six cubes. A dataset of expert demonstrations is available for the training task and the model is evaluated on its ability to predict the block positions and rank action sequences in expert demonstrations of the zero-shot transfer tasks. The block position prediction metric measures if ground-truth information about the positions of all blocks can be decoded from the learned latent space. The action ranking metric measures if the model can distinguish action sequences that successfully solve a task from perturbed action sequences that do not. See Appendix A.1 for additional information.

Result: Factored world models with soft action attention outperform the baseline models when the number of GNN layers is set to one or two (Figure 1). As the number of layers increases to three and four, the performance of the two models converges. The attention weights predicted by Soft Attention for a model with one graph neural network layer are shown in Table 2 in the Appendix. The model correctly identifies which object is being manipulated. Additionally, the model pays attention to where a particular cube is being placed. When building a tower, the cubes below the placed object also receive the information about the `place` action with a low (but non-zero) weight. We believe this behavior emerged because placement of a cube can often move the cubes below.

Furthermore, we hypothesize that the diminishing benefit of soft attention as we add more GNN layers can be explained by the increasing capacity of the model. A single GNN layer can benefit from the pre-computed information about which object will be affected. On the other hand, a stack of four GNN layers receives the information about the action at each layer, making it easier to distinguish affected and unaffected objects within the transition model.

5 CONCLUSION

We propose two attention mechanisms for binding actions to objects, soft attention and hard attention, and validate them on the task of learning structured world models in five environments. As one of our reviewers pointed out, “there seems to be a dilemma between soft and hard attention – the soft attention is more expressive, however, it also under-performs hard attention in some experiments.” We view this trade-off from the perspective of expressiveness versus bias. Hard attention provides a strong bias (that all actions move exactly one object), which allows C-SWM to converge to the right solution in the grid-world environments. Conversely, soft attention is more flexible in providing information about the action to several objects. This construction fits our simulated robotic manipulation task, where several cubes might move at once when they are being stacked.

ACKNOWLEDGEMENTS

This work is supported in part by NSF 1724257, NSF 1724191, NSF 1763878, NSF 1750649, NSF 1835309, NSF 2107256 and NASA 80NSSC19K1474.

REFERENCES

- Victor Bapst, Alvaro Sanchez-Gonzalez, Carl Doersch, Kimberly L. Stachenfeld, Pushmeet Kohli, Peter W. Battaglia, and Jessica B. Hamrick. Structured agents for physical construction. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 464–474. PMLR, 2019.
- Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinícius Flores Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Çağlar Gülçehre, H. Francis Song, Andrew J. Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani, Kelsey R. Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matthew Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. *CoRR*, abs/1806.01261, 2018.
- Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *J. Artif. Intell. Res.*, 47:253–279, 2013.
- Ondrej Biza, Elise van der Pol, and Thomas Kipf. The impact of negative sampling on contrastive structured world models. *ICML 2021 Workshop: Self-Supervised Learning for Reasoning and Perception*, 2021.
- Ondrej Biza, Thomas Kipf, David Klee, Robert Platt, Jan-Willem van de Meent, and Lawson L. S. Wong. Factored world models for zero-shot generalization in robotic manipulation, 2022.
- M. Gori, G. Monfardini, and F. Scarselli. A new model for learning in graph domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 2, pp. 729–734 vol. 2, 2005.
- Anirudh Goyal, Aniket Rajiv Didolkar, Nan Rosemary Ke, Charles Blundell, Philippe Beaudoin, Nicolas Heess, Michael Curtis Mozer, and Yoshua Bengio. Neural production systems. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021.
- Klaus Greff, Sjoerd van Steenkiste, and Jürgen Schmidhuber. On the binding problem in artificial neural networks. *CoRR*, abs/2012.05208, 2020.
- William H Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codel, Manuela Veloso, and Ruslan Salakhutdinov. Minerl: A large-scale dataset of minecraft demonstrations. *arXiv preprint arXiv:1907.13440*, 2019.
- Michael Janner, Sergey Levine, William T. Freeman, Joshua B. Tenenbaum, Chelsea Finn, and Jiajun Wu. Reasoning about physical interactions with object-oriented prediction and planning. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- Thomas N. Kipf, Elise van der Pol, and Max Welling. Contrastive learning of structured world models. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- Richard Li, Allan Jabri, Trevor Darrell, and Pulkit Agrawal. Towards practical multi-object manipulation using relational reinforcement learning. In *2020 IEEE International Conference on Robotics and Automation, ICRA 2020, Paris, France, May 31 - August 31, 2020*, pp. 4051–4058. IEEE, 2020.
- Marsha Lovett and John R. Anderson. Thinking as a production system. *The Cambridge handbook of thinking and reasoning*, 2005.

Shikhar Sharma, Ryan Kiros, and Ruslan Salakhutdinov. Action recognition using visual attention. *CoRR*, abs/1511.04119, 2015.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008, 2017.

Rishi Veerapaneni, John D. Co-Reyes, Michael Chang, Michael Janner, Chelsea Finn, Jiajun Wu, Joshua B. Tenenbaum, and Sergey Levine. Entity abstraction in visual model-based reinforcement learning. In Leslie Pack Kaelbling, Danica Kragic, and Komei Sugiura (eds.), *3rd Annual Conference on Robot Learning, CoRL 2019, Osaka, Japan, October 30 - November 1, 2019, Proceedings*, volume 100 of *Proceedings of Machine Learning Research*, pp. 1439–1456. PMLR, 2019.

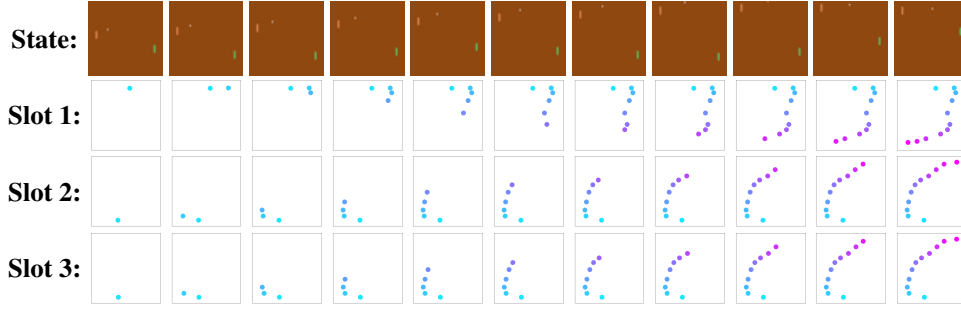


Figure 2: Visualization of encoded states for a single episode of Atari Pong. Each state is encoded into a latent state with three slots, and we color encodings from blue to purple as the time progresses.

A APPENDIX

A.1 EVALUATION METRICS

C-SWM Hits@1 and Mean Reciprocal Rank: We follow the evaluation setup from [Kipf et al. \(2020\)](#). An evaluation dataset of episodes of length 10 is collected. The agent’s encoder is used to encode all states in the evaluation dataset into latent states. Then the agent predicts the next latent state $t = 1$, $t = 5$ or $t = 10$ steps into the future given a sequence of 1, 5 or 10 actions respectively. The predicted latent state is then compared to the actual encoded state at timestep t (the positive example) and to all states at timestep t from other episodes (the negative example). The Hits@1 metric simply counts the fraction of times the predicted latent state was closer to the positive example than to all other negative examples. The Mean Reciprocal Rank (MRR) metric computes

$$\frac{1}{|D|} \sum_{n=1}^{|D|} \frac{1}{\text{rank}_n}, \quad (7)$$

where D is the evaluation dataset and rank_n is the index of the positive examples in the sorted list of distances between the predicted latent state and the set of the positive and all negative examples.

The original evaluation code does not always handle duplicate states in different episodes well. It is assumed that two identical images will be encoded into the same latent state (i.e. all float32 values will be exactly equal), but this does not happen in practice even when pytorch is set to a deterministic mode. Duplicate states are infrequent in 2D Shapes and 3D Cubes, but fairly common in Atari. In the Atari experiments, we check for duplicate states and remove negative states that are identical to a particular positive state.

Slot Correlation: we compute the absolute value of the Pearson correlation coefficient between every pair of object slots. This value is averaged over all pairs of object slots and over all states in the evaluation set. The resulting value is bounded between zero and one, where zero means no correlation between slots and one means that all slots capture the same information.

FWM Block Position RMSE and Action Sequence Hits@1: We use the same metrics as [Biza et al. \(2022\)](#). To the block position prediction, we train a Multi-Layer Perceptron to predict the spatial (x, y, z) position of each block from the latent space of FWMs. Only the MLP is trained, all other weights are frozen. Then, we compute the root mean squared error between the actual and predicted block positions. The dataset used in this evaluation is the dataset of expert demonstrations for the four zero-shot transfer tasks.

In action sequence ranking, we generate one correct and ten incorrect trajectories for each zero-shot transfer task. We also provide the model with an example of a goal state. The model encodes the goal state and predicts the resulting latent state of all eleven action sequences. If the predicted result for the correct trajectory is the closest to the encoded goal state, the prediction of the model is considered correct. Hits@1 once again measures the fraction of correct predictions.

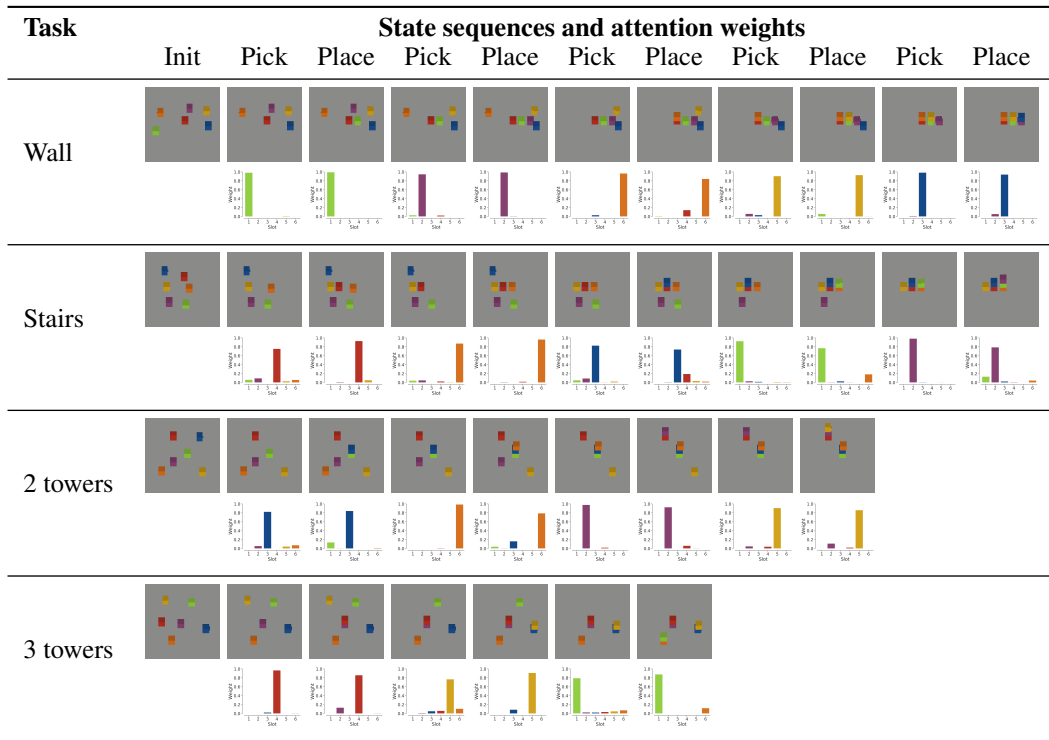


Table 2: Visualization of attention weights for the four zero-shot transfer tasks in our pick-and-place experiment. Individual bars are connected to objects in the images by their colors.