
Grounding Lifted PDDL Action Models

Masataro Asai

IBM Research, MIT-IBM Watson AI Lab
Cambridge, USA

Abstract

We propose FOSAE++, an unsupervised end-to-end neural system that generates a compact discrete state transition model (dynamics / action model) from raw visual observations. Our representation can be exported to Planning Domain Description Language (PDDL), allowing symbolic state-of-the-art classical planners to perform high-level task planning. FOSAE++ expresses states and actions in First-Order Logic (FOL). It is the first unsupervised neural system that fully supports FOL in PDDL action modeling, while existing systems are limited to continuous, propositional, or property-based representations, and/or require labeled actions.

1 Introduction

Recently, Asai and Fukunaga [2018, Latplan] has opened the door to applying symbolic Classical Planning systems to a wide variety of raw, noisy data. Latplan uses discrete variational autoencoders to generate propositional latent states and its dynamics (action model) directly from images. Unlike existing work, which requires several machine learning pipelines (SVM/decision trees) and labeled inputs (e.g., a sequence of high-level options) [Konidaris et al., 2014], Latplan is an end-to-end unsupervised neural network that requires no manually labeled inputs.

In this paper, we obtain a *lifted action model* expressed in First-Order Logic (FOL). In propositional action models, the representation is a fixed-sized binary array and does not transfer to a different or a dynamically changing environment with a varying number of objects. In contrast, lifted FOL representations are generalized over objects and environments, as we demonstrate in Blocksworld with different number of blocks, or Sokoban with different map sizes. We propose *Lifted First-Order Space AutoEncoder* (FOSAE++) neuro-symbolic architecture, which learns a lifted PDDL action model by integrating and extending (1) the First-Order State AutoEncoder [Asai, 2019, FOSAE] that extends Discrete VAE to generate predicate symbols, (2) the Cube-Space AutoEncoder [Asai and Muise, 2020, CSAE] which regularizes the latent space to be directly exportable to a learned propositional PDDL model [Fikes et al., 1972], (3) and the Neural Logic Machine [Dong et al., 2019, NLM] which can process FOL statements. The overall task of our system is illustrated in Fig. 1. The system takes a *transition dataset* containing a set of pairs of raw observations which are single time step away. Each observation consists of multiple visual segmentations of the objects. It learns a lifted action model of the environment by generating propositional/predicate/action symbols and emits a PDDL [Haslum et al., 2019] encoding for state-of-the-art planning systems. Table 1 contains a taxonomy of existing model acquisition systems in chronological order. FOSAE++ is the first system that satisfies all features readily available in symbolic action model acquisition systems, while not relying on human-derived symbols.

2 Preliminaries and Background

We denote a concatenation of tensors \mathbf{a} and \mathbf{b} in the last axis by $\mathbf{a}; \mathbf{b}$, and the i -th data point of a dataset by a superscript i which we may omit for clarity. **We use the same symbol for a set and its**

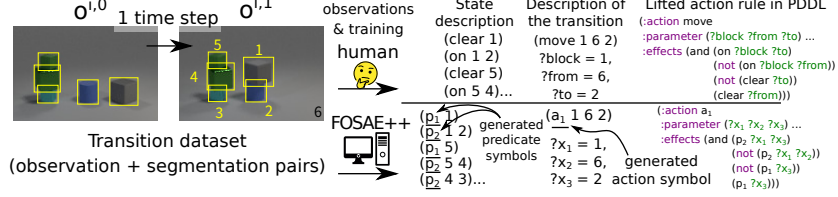


Figure 1: FOSAE++’s learning task compared to the action modeling by humans. All symbols are anonymous (GENSYM’d) due to the unsupervised nature.

	Symbol Generation			Compatibility			
	Propo- sitional	Multi-arity predicates	Action	End-to- End NN	Search algorithms	Propositional PDDL	Lifted PDDL
Symbolic approaches ¹	m	m	m	no	yes	yes	yes
[Konidaris et al., 2014]	yes	u	m ²	no	yes	yes	no
Latplan [Asai and Fukunaga, 2018]	yes	u	yes	yes	yes	no	no
CausalInfoGAN [Kurutach et al., 2018]	yes	u	u	yes	i ²	no	no
FOSAE [Asai, 2019]	yes	yes	m	yes	yes	no	no
[James et al., 2020b]	yes	i ²	m ²	no	yes	yes	i ²
CSAE [Asai and Muise, 2020]	yes	u	yes	yes	yes	yes	no
FOSAE++	yes	yes	yes	yes	yes	yes	yes

Table 1: Taxonomy of action model acquisition systems in chronological order. m:manual, u:unused, i:incomplete.

size (e.g., S , and not $|S|$) to avoid the clutter. $\mathbb{B} = [0, 1] \subset \mathbb{R}$. We assume background knowledge of discrete VAEs with continuous relaxations, such as Gumbel-Softmax (GS) and Binary-Concrete (BC) [Jang et al., 2017, Maddison et al., 2017]. Their activations are denoted as GS and BC, respectively. We also assume the basic knowledge of classical planning available in standard AI textbooks Russell et al. [1995]. For those backgrounds, please refer to the Appendix Sec. B.

Latplan [Asai and Fukunaga, 2018] learns a propositional state representation and transition rules entirely from image-based observations of the environment with discrete VAEs and solves the problem using a classical planner. Latplan is trained on a set of raw state transition observations randomly sampled from the environment. The i -th data $(o^{i,0}, o^{i,1})$ is a pair of observations made before and after an unknown high-level action is performed. Latplan then processes a *planning input* (o^I, o^G) , a pair of raw images corresponding to an initial and goal state of the environment. The output of Latplan is a data sequence representing the plan execution $(o^I, \dots o^G)$ that reaches o^G from o^I .

Latplan works in 4 steps. **(1)** A Binary-Concrete VAE called *State AutoEncoder* (SAE) learns to represent the input data o as a propositional/boolean states z . The encoder then generates propositional transitions $\{\dots (z^{i,0}, z^{i,1}) \dots\}$ from $\{\dots (o^{i,0}, o^{i,1}) \dots\}$. **(2)** A Gumbel-Softmax VAE with a skip connection called *Action AutoEncoder* (AAE) generates action symbols and learns transition rules (action model/dynamics) from the propositional transitions. **(3)** It generates a propositional classical planning problem (z^I, z^G) from a planning input (o^I, o^G) and export it into PDDL files together with the action model obtained in (2), which are then solved by Fast Downward [Helmert, 2006] classical planning solver which returns a plan. **(4)** Finally, the plan’s intermediate states are visualized into a step-by-step, human-comprehensible images by the SAE decoder.

AAE performs clustering on the state transitions with the maximum number of clusters A specified as a hyperparameter. The cluster ID is used as its action symbol. The clustering is performed by its encoder, $\text{ACTION}(z^{i,0}, z^{i,1}) = a^i$, which takes a propositional state pair $(z^{i,0}, z^{i,1})$ and returns a one-hot vector a^i of A categories using Gumbel-Softmax. AAE’s decoder takes the current state $z^{i,0}$ and the action a^i and reconstructs a successor state $\tilde{z}^{i,1} \approx z^{i,1}$, acting as a progression/dynamics

¹[Yang et al., 2007, Cresswell et al., 2013, Aineto et al., 2018, Zhuo et al., 2019, Cresswell and Gregory, 2011, Mourão et al., 2012, Zhuo and Kambhampati, 2013]

²Konidaris et al. [2014] requires sequences of high-level options to learn from, such as [move, move, interact, ...] in Playroom domain. Causal InfoGAN cannot deterministically enumerate all successors (a requirement for search completeness) due to the lack of action symbols and should sample the successors. James et al. [2020b]’s PDDL output is limited to unary predicates / properties of objects, thus cannot model the interactions between objects. Also, it requires sequences of high-level options such as [WalkToltem, AttachBlock, WalkToltem, ...] in the Minecraft domain.

$\text{APPLY}(\mathbf{a}^i, \mathbf{z}^{i,0}) = \tilde{\mathbf{z}}^{i,1}$. AAE is trained to minimize the successor reconstruction loss $\ell(\tilde{\mathbf{z}}^{i,1}, \mathbf{z}^{i,1})$. Asai and Muise [2020] proposed to combine SAE and AAE into a single network called Space AE, which additionally reconstructs a predicted successor $\tilde{\mathbf{z}}^{i,1}$ to $\tilde{\tilde{\mathbf{z}}}^{i,1}$.

APPLY in AAE/Space AE is an arbitrary neural black-box function that does not directly translates to PDDL, preventing efficient search with PDDL-based solvers. Cube-Space AE [Asai and Muise, 2020] addresses this issue by Back-To-Logit technique (BTL) which restricts the dynamics to a Karman Filter and discretize it with BC ($\tilde{\mathbf{z}}^{i,1} = \text{BC}(\mathbf{F}\mathbf{z}^{i,0} + \mathbf{B}\mathbf{a}^i + \mathbf{w})$), and directly compile it into STRIPS/PDDL language. See Appendix Sec. C for detailed discussion.

First-Order State AutoEncoder [Asai, 2019, FOSAE] is an autoencoder that represents a set of *object feature vectors* in a latent predicate representation. Unlike prior work on relational modeling [Santoro et al., 2017, Zambaldi et al., 2019, Battaglia et al., 2018], this system obtains discrete logical predicates compatible with symbolic systems. The input is similar to the setting of Ugur and Piater [2015] and James et al. [2020b]: Each object feature vector represents each object in an arbitrary complex manner. In this paper, we use segmented&resized pixels and bounding box information (x, y and width, height). Let $\mathbf{o}_n \in \mathbb{R}^F$ be an F -dimensional feature vector representing each object and $\mathbf{o} = (\mathbf{o}_1, \dots, \mathbf{o}_O) \in \mathbb{R}^{O \times F}$ be the input matrix representing O objects. FOSAE can generate a *multi-arity latent representation*. We denote a set of predicates of arity n as P/n (Prolog notation) and its propositions as $P/n(O)$. We denote a binary tensor representation of $P/n(O)$ as $\mathbf{z}/n \in \mathbb{B}^{O \times \dots \times O \times P/n}$, and the latent space is a tuple $\mathbf{z} = (\mathbf{z}/1, \dots, \mathbf{z}/N)$, where N is the largest arity. To convert \mathbf{o} to n -ary predicates \mathbf{z}/n (e.g., $n = 2$ for relationships), we apply a 1D pointwise convolutional filter f_n of P/n output features over each concatenated objects, e.g., $\mathbf{z}/2_{ij} = \text{BC}(f_2(\mathbf{o}_i; \mathbf{o}_j)) \in \mathbb{B}^{P/2}$. The latent tensors $(\mathbf{z}/1, \dots, \mathbf{z}/N)$ can be flattened, concatenated, decoded to the reconstruction $\tilde{\mathbf{o}} \in \mathbb{R}^{O \times F}$ by any decoder network (MLP in the original paper).

The NLM [Dong et al., 2019] is a neural Inductive Logic Programming (ILP) [Muggleton, 1991] system whose inputs and outputs are hand-crafted binary tensors representing propositional groundings of FOL statements. These binary tensors are in the same multi-arity representation $(\mathbf{z}/1, \dots, \mathbf{z}/N)$ used by FOSAE. An NLM output is denoted as $\text{NLM}_{Q,\sigma}(\mathbf{z}) = (\text{COMPOSE}_{Q,\sigma}(\mathbf{z}, 1), \dots, \text{COMPOSE}_{Q,\sigma}(\mathbf{z}, N))$, where Q is the number (hyperparameter) of output features, and σ is a nonlinearity. See Appendix Sec. B.3 for details.

3 Lifting the Action Model

In order to obtain a lifted action model in an unsupervised setting, there are three requirements: (1) white-box action model which is trivially convertible to STRIPS formalism, (2) invariance to the number/order of objects, (3) unsupervised generation of multi-arity predicate symbols. To our knowledge, no existing systems fulfill all requirements: FOSAE lacks 1 and 2, CSAE lacks 2 and 3, and NLM lacks 1 and 3 (designed for hand-crafted symbolic data). Our FOSAE++ addresses all issues at once. Its overall design follows the CSAE, which ENCODE s the input, identifies the action \mathbf{a}^i , APPLY-es the action, then DECODE s the results. It uses FOSAE’s encoder to generate multi-arity latent representation, which is then consumed by NLMs used as a basic building block of other components. *This architecture’s key element is a new component PARAMS and a unique pair of operations called BIND and UNBIND*, which intuitively reflects the structure of lifted action models. Suppose we model a lifted action (move ?block ?from ?to) in Blocksworld (Fig. 1), with effects such as (on ?block ?to). Since a lifted model always refers to the objects through its parameters such as ?to, it cannot affect the objects that does not appear in the parameter list, e.g., action (move a b c) cannot affect (on d c) because $d \notin \{a, b, c\}$. *We represent this restriction as differentiable operations.*

PARAMS is an attention network $\mathbf{x} = \text{PARAMS}(\mathbf{z}^{i,0}; \mathbf{z}^{i,1}) = (\mathbf{x}_1, \dots, \mathbf{x}_{\#a})$ which attends to objects, essentially learning the parameters of the action. $\#a$ is a hyperparameter. Each parameter \mathbf{x}_i is a one-hot hard attention vector activated by Gumbel-Softmax, thus can later extract a specific object from the object feature matrix \mathbf{o} . PARAMS consists of NLM layers, ending with $\text{NLM}_{\#a, \text{IDENTITY}}$. We then extract its unary part of the results ($\rightarrow \mathbb{R}^{O \times \#a}$), transpose it ($\rightarrow \mathbb{R}^{\#a \times O}$), then apply $\#a$ -way Gumbel-Softmax of O categories ($\rightarrow \mathbb{B}^{\#a \times O}$). As a result, the output attends to $\#a$ objects in total.

\mathbf{x} is used by unique operations named BIND and UNBIND. Recall that the predicates in the effects can refer only to the specified action parameters. Therefore, we *limit the dynamics to the objects attended by \mathbf{x}* . We iteratively soft-extracts the sub-axes of \mathbf{z}/n attended by \mathbf{x} using matrix operation

References

- Diego Aineto, Sergio Jiménez, and Eva Onaindia. Learning STRIPS Action Models with Classical Planning. In *Proc. of the International Conference on Automated Planning and Scheduling(ICAPS)*, 2018.
- Vidal Alcázar, Daniel Borrajo, Susana Fernández, and Raquel Fuentetaja. Revisiting Regression in Planning. In *Proc. of International Joint Conference on Artificial Intelligence (IJCAI)*, 2013.
- Garrett Andersen and George Konidaris. Active Exploration for Learning Symbolic Representations. In *Advances in Neural Information Processing Systems*, pages 5009–5019, 2017.
- Masataro Asai. Photo-Realistic Blocksworld Dataset. *arXiv:1812.01818*, 2018. URL <http://arxiv.org/abs/1812.01818>.
- Masataro Asai. Unsupervised Grounding of Plannable First-Order Logic Representation from Images. In *Proc. of the International Conference on Automated Planning and Scheduling(ICAPS)*, July 2019.
- Masataro Asai and Alex Fukunaga. Classical Planning in Deep Latent Space: Bridging the Subsymbolic-Symbolic Boundary. In *Proc. of AAAI Conference on Artificial Intelligence*, 2018. URL <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16302>.
- Masataro Asai and Hiroshi Kajino. Towards Stable Symbol Grounding with Zero-Suppressed State AutoEncoder. In *Proc. of the International Conference on Automated Planning and Scheduling(ICAPS)*, July 2019.
- Masataro Asai and Christian Muise. Learning Neural-Symbolic Descriptive Planning Models via Cube-Space Priors: The Voyage Home (to STRIPS). In *Proc. of International Joint Conference on Artificial Intelligence (IJCAI)*, 2020. URL <https://arxiv.org/abs/2004.12850>.
- Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- Stephen Cresswell and Peter Gregory. Generalised Domain Model Acquisition from Action Traces. In *Proc. of the International Conference on Automated Planning and Scheduling(ICAPS)*, 2011. URL <http://aaai.org/ocs/index.php/ICAPS/ICAPS11/paper/view/2712>.
- Stephen Cresswell, Thomas Leo McCluskey, and Margaret Mary West. Acquiring planning domain models using *LOCM*. *Knowledge Eng. Review*, 28(2):195–213, 2013. doi: 10.1017/S0269888912000422.
- Joseph Culberson. Sokoban is PSPACE-complete. In *Proceedings in Informatics 4, International Conference on Fun with Algorithms*, pages 65–76. Carleton Scientific, June 1998.
- Honghua Dong, Jiayuan Mao, Tian Lin, Chong Wang, Lihong Li, and Denny Zhou. Neural Logic Machines. In *Proc. of the International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=B1xY-hRctX>.
- Martin Engelcke, Adam R Kosiorek, Oiwi Parker Jones, and Ingmar Posner. GENESIS: Generative Scene Inference and Sampling with Object-Centric Latent Representations. In *Proc. of the International Conference on Learning Representations*, 2020.
- SM Ali Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, David Szepesvari, Geoffrey E Hinton, et al. Attend, Infer, Repeat: Fast Scene Understanding with Generative Models. In *Advances in Neural Information Processing Systems*, pages 3225–3233, 2016.
- Richard E Fikes, Peter E. Hart, and Nils J. Nilsson. Learning and Executing Generalized Robot Plans. *Artificial Intelligence*, 3(1-3):251–288, 1972. doi: 10.1016/0004-3702(72)90051-3. URL [http://dx.doi.org/10.1016/0004-3702\(72\)90051-3](http://dx.doi.org/10.1016/0004-3702(72)90051-3).
- Kunihiko Fukushima. Neocognitron: A Self-Organizing Neural Network Model for a Mechanism of Pattern Recognition Unaffected by Shift in Position. *Biological Cybernetics*, 36(4), 1980.

- Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Christopher Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-Object Representation Learning with Iterative Variational Inference. In *Proc. of the International Conference on Machine Learning*, pages 2424–2433, 2019.
- Patrik Haslum, Nir Lipovetzky, Daniele Magazzeni, and Christian Muise. *An Introduction to the Planning Domain Definition Language*, volume 13. Morgan & Claypool Publishers, 2019.
- Malte Helmert. The Fast Downward Planning System. *J. Artif. Intell. Res.(JAIR)*, 26:191–246, 2006. URL <http://www.aaai.org/Papers/JAIR/Vol26/JAIR-2606.pdf>.
- Malte Helmert and Carmel Domshlak. Landmarks, Critical Paths and Abstractions: What’s the Difference Anyway? In *Proc. of the International Conference on Automated Planning and Scheduling(ICAPS)*, 2009. URL <http://www.aaai.org/ocs/index.php/ICAPS/ICAPS09/paper/viewPDFInterstitial/735/1107>.
- De-An Huang, Danfei Xu, Yuke Zhu, Animesh Garg, Silvio Savarese, Li Fei-Fei, and Juan Carlos Niebles. Continuous Relaxation of Symbolic Planner for One-Shot Imitation Learning. In *Proc. of IEEE International Workshop on Intelligent Robots and Systems (IROS)*, pages 2635–2642. IEEE, 2019.
- Sergey Ioffe and Christian Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *Proc. of the International Conference on Machine Learning*, pages 448–456, 2015. URL <http://jmlr.org/proceedings/papers/v37/ioffe15.html>.
- Steven James, Benjamin Rosman, and George Konidaris. Learning Portable Representations for High-Level Planning. In *Proc. of the International Conference on Machine Learning*, 2020a.
- Steven James, Benjamin Rosman, and George Konidaris. Learning Object-Centric Representations for High-Level Planning in Minecraft. In *Proceedings of Workshop on Object-Oriented Learning at ICML*, 2020b.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical Reparameterization with Gumbel-Softmax. In *Proc. of the International Conference on Learning Representations*, 2017.
- Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1988–1997. IEEE, 2017.
- Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960.
- Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. In *Proc. of the International Conference on Learning Representations*, 2013.
- George Konidaris, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. Constructing Symbolic Representations for High-Level Planning. In *Proc. of AAAI Conference on Artificial Intelligence*, pages 1932–1938, 2014. URL <http://www.aaai.org/ocs/index.php/AAAI/AAAI14/paper/view/8424>.
- George Konidaris, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. Symbol Acquisition for Probabilistic High-Level Planning. In *Proc. of International Joint Conference on Artificial Intelligence (IJCAI)*, pages 3619–3627, 2015. URL <http://ijcai.org/Abstract/15/509>.
- George Konidaris, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. From Skills to Symbols: Learning Symbolic Representations for Abstract High-Level Planning. *J. Artif. Intell. Res.(JAIR)*, 61: 215–289, 2018. doi: 10.1613/jair.5575. URL <https://doi.org/10.1613/jair.5575>.
- Thanard Kurutach, Aviv Tamar, Ge Yang, Stuart Russell, and Pieter Abbeel. Learning Plannable Representations with Causal InfoGAN. In *Advances in Neural Information Processing Systems*, 2018.
- Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables. In *Proc. of the International Conference on Learning Representations*, 2017.

- Kira Mourão, Luke S. Zettlemoyer, Ronald P. A. Petrick, and Mark Steedman. Learning STRIPS Operators from Noisy and Incomplete Observations. In *Proc. of the International Conference on Uncertainty in Artificial Intelligence*, pages 614–623, 2012. URL https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=2322&proceeding_id=28.
- Stephen Muggleton. Inductive Logic Programming. *New generation computing*, 8(4):295–318, 1991.
- Vinod Nair and Geoffrey E Hinton. Rectified Linear Units Improve Restricted Boltzmann Machines. In *Proc. of the International Conference on Machine Learning*, 2010.
- Charles Payan. On the Chromatic Number of Cube-Like Graphs. *Discrete mathematics*, 103(3), 1992.
- Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- Raymond Reiter. On Closed World Data Bases. In *Readings in Artificial Intelligence*, pages 119–140. Elsevier, 1981.
- Stuart Jonathan Russell, Peter Norvig, John F Canny, Jitendra M Malik, and Douglas D Edwards. *Artificial Intelligence: A Modern Approach*, volume 2. Prentice hall Englewood Cliffs, 1995.
- Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Tim Lillicrap. A Simple Neural Network Module for Relational Reasoning. In *Advances in Neural Information Processing Systems*, pages 4967–4976, 2017.
- Tom Silver and Rohan Chitnis. PDDL Gym: Gym Environments from PDDL Problems, 2020.
- Emre Ugur and Justus Piater. Bottom-up Learning of Object Categories, Action Effects and Logical Rules: From Continuous Manipulative Exploration to Symbolic Planning. In *Proc. of IEEE International Conference on Robotics and Automaton (ICRA)*, pages 2627–2633. IEEE, 2015.
- Arash Vahdat, Evgeny Andriyash, and William Macready. DVAE#: Discrete variational autoencoders with relaxed Boltzmann priors. In *Advances in Neural Information Processing Systems*, 2018a.
- Arash Vahdat, William G Macready, Zhengbing Bian, Amir Khoshaman, and Evgeny Andriyash. DVAE++: Discrete variational autoencoders with overlapping transformations. *arXiv:1802.04920*, 2018b.
- Aaron van den Oord, Oriol Vinyals, et al. Neural Discrete Representation Learning. In *Advances in Neural Information Processing Systems*, 2017.
- Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric Cross Entropy for Robust Learning with Noisy Labels. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition*, pages 322–330, 2019.
- Kai Xu, Akash Srivastava, and Charles Sutton. Variational Russian Roulette for Deep Bayesian Nonparametrics. In *International Conference on Machine Learning*, pages 6963–6972, 2019.
- Qiang Yang, Kangheng Wu, and Yunfei Jiang. Learning Action Models from Plan Examples using Weighted MAX-SAT. *Artificial Intelligence*, 171(2-3):107–143, 2007. doi: 10.1016/j.artint.2006.11.005. URL <http://dx.doi.org/10.1016/j.artint.2006.11.005>; <http://dblp.uni-trier.de/rec/bib/journals/ai/YangWJ07>.
- Vinicius Zambaldi, David Raposo, Adam Santoro, Victor Bapst, Yujia Li, Igor Babuschkin, Karl Tuyls, David Reichert, Timothy Lillicrap, Edward Lockhart, et al. Deep Reinforcement Learning with Relational Inductive Biases. In *Proc. of the International Conference on Learning Representations*, 2019.
- Hankz Hankui Zhuo and Subbarao Kambhampati. Action-Model Acquisition from Noisy Plan Traces. In *Proc. of International Joint Conference on Artificial Intelligence (IJCAI)*, 2013.
- Hankz Hankui Zhuo, Jing Peng, and Subbarao Kambhampati. Learning Action Models from Disordered and Noisy Plan Traces. *arXiv:1908.09800*, 2019.

A Related Work

James et al. [2020a,b] build on existing work [Konidaris et al., 2014, 2015, 2018] to find an *ego-centric* logical representation that is invariant to the state of the observer in a complex environment, or an *object-centric*, property-based logical representation. Both representations are special cases of First-Order Logic representation where all predicates are unary. These approaches also require a training input that already contains action symbols (e.g., toggle-door). Andersen and Konidaris [2017] obtains an MDP-based PPDDL (Probabilistic PDDL) model using Active Learning. Incorporating Active Learning in FOSAE++ for self data collection is future work.

Huang et al. [2019] reported that direct discretization approach performed worse than an approach that plans in the probabilistic representation obtained by neural networks. However, the question of continuous versus discrete is not conclusive. They used a naive rounding-based discretization which may cause a significant amount of information loss compared to the state-of-the-art discrete variational autoencoders. Furthermore, they rely on mappings from raw inputs to hand-coded symbolic representations that require supervised training.

Object-based input we used can be obtained from state-of-the-art object-recognition methods such as YOLO [Redmon and Farhadi, 2018]. More recent systems [Greff et al., 2019, Engelcke et al., 2020] can handle shapes not limited to rectangles, and can be trained unsupervised. Connecting our network with these systems is an exciting avenue of future work.

Eslami et al. [2016] and Xu et al. [2019] proposed an attention-based and a Bayesian approach to the problem of variable-sized object reconstruction. Their work do not address the state dynamics and is orthogonal to our work. The primary difference from our approach is that they studied on sequentially storing and retrieving a variable amount of information into/from a fixed-sized array, while we store them in a latent representation of the corresponding size.

B Extended Backgrounds

B.1 Lifted STRIPS/PDDL Planning

Planning Domain Description Language (PDDL) is a modeling language for a Lifted STRIPS planning formalism [Fikes et al., 1972] and its extensions [Haslum et al., 2019]. Let $\mathcal{F}(T)$ be a formula consisting of logical operations $\{\wedge, \neg\}$ and a set of terms T . For example, when $T = \{have(I, food), full(I)\}$, then $have(I, food) \wedge \neg full(I) \in \mathcal{F}(T)$. We denote a lifted STRIPS planning problem as a 5-tuple $\langle O, P, A, I, G \rangle$. O is a set of objects ($\ni food$), P is a set of predicates ($\ni full(x)$), and A is a set of lifted actions ($\ni eat$). Each predicate $p \in P$ has an arity $\#p \geq 0$. Predicates are instantiated/grounded into propositions $P(O) = \bigcup_{p \in P} (\{p\} \times O^{\#p} \times O)$, such as $have(I, food)$. A state $s \subseteq P(O)$ represents truth assignments to the propositions, e.g., $s = \{have(I, food)\}$ represents $have(I, food) = \top$. We can also represent it as a bitvector of size $\sum_p O^{\#p}$.

Each lifted action $a(X) \in A$ has an arity $\#a$ and parameters $X = (x_1, \dots, x_{\#a})$, such as $eat(x_1, x_2)$. Lifted actions are instantiated into *ground actions* $A(O) = \bigcup_{a \in A} (\{a\} \times O^{\#a} \times O)$, such as $eat(I, food)$. $a(X)$ is a 3-tuple $\langle PRE(a), ADD(a), DEL(a) \rangle$, where $PRE(a), ADD(a), DEL(a) \in \mathcal{F}(P(X))$ are preconditions, add-effects, and delete-effects: e.g., $eat(x_1, x_2) = \langle \{have(x_1, x_2)\}, \{full(x_1)\}, \{have(x_1, x_2)\} \rangle$. The semantics of these three elements are as follows: A ground action $a^\dagger \in A(O)$ is *applicable* when a state s satisfies $PRE(a^\dagger)$, i.e., $PRE(a^\dagger) \subseteq s$, and applying an action a^\dagger to s yields a new successor state $a^\dagger(s) = (s \setminus DEL(a^\dagger)) \cup ADD(a^\dagger)$, e.g., $eat(I, food) = \text{“I can eat a food when I have one, and if I eat one I am full but the food is gone.”}$ Finally, $I, G \subseteq P(O)$ are the initial state and a goal condition, respectively. The task of classical planning is to find a *plan* $(a_1^\dagger, \dots, a_n^\dagger)$ which satisfies $a_n^\dagger \circ \dots \circ a_1^\dagger(I) \subseteq G$.

B.2 Discrete Variational Autoencoders

Variational AutoEncoder (VAE) is a framework for reconstructing an observation x from a compact latent representation z that follows a certain prior distribution. Training is performed by maximiz-

ing the sum of the reconstruction loss and the KL divergence between the latent random distribution $q(z|x)$ and the target distribution $p(z)$, which gives a lower bound for the likelihood $p(x)$ [Kingma and Welling, 2013]. While typically $p(z)$ is a Normal distribution $\mathcal{N}(0, 1)$, Gumbel-Softmax (GS) VAE [Jang et al., 2017] and Binary-Concrete (BC) VAE [Maddison et al., 2017] use a discrete, uniform categorical/Bernoulli(0.5) distribution, and further approximate it with a continuous relaxation by introducing a temperature parameter τ that is annealed down to 0.

We denote corresponding activation function in the latent space as $\text{GS}(\mathbf{l})$ and $\text{BC}(l)$, where \mathbf{l} and l each represents a logit vector and scalar. A latent vector $\mathbf{z} \in [0, 1]^C$ of GS-VAE is computed from a logit vector $\mathbf{l} \in \mathbb{R}^C$ by $\mathbf{z} = \text{GS}(\mathbf{l}) = \text{SOFTMAX}((\mathbf{l} + \text{GUMBEL}(0, 1))/\tau)$, where C is the number of categories, $\text{GUMBEL}(0, 1) = -\log(-\log u)$ and $u \sim \text{UNIFORM}(0, 1) \in [0, 1]^C$.

A latent scalar z of BC-VAE is computed from a logit scalar $l \in \mathbb{R}$ by $z = \text{BC}(l) = \text{SIGMOID}((l + \text{LOGISTIC}(0, 1))/\tau)$, where $\text{LOGISTIC}(0, 1) = \log u - \log(1 - u)$ and $u \sim \text{UNIFORM}(0, 1) \in \mathbb{R}$.

Both functions converge to discrete functions at the limit $\tau \rightarrow 0$: $\text{GS}(\mathbf{l}) \rightarrow \arg \max(\mathbf{l})$ (in one-hot representation) and $\text{BC}(l) \rightarrow \text{STEP}(l) = (l < 0) ? 0 : 1$. Typically, GS-VAE and BC-VAE contains multiple latent vectors / latent scalars to model the complex input.

There are more complex variations such as VQVAE van den Oord et al. [2017], DVAE++Vahdat et al. [2018b], DVAE# Vahdat et al. [2018a]. They may contribute to the stable performance, but we leave the task of faster / easier training to the future work.

B.3 Neural Logic Machine

The NLM [Dong et al., 2019] is a neural Inductive Logic Programming (ILP) [Muggleton, 1991] system based on First-Order Logic and the Closed-World Assumption [Reiter, 1981]. Given a set of base predicates grounded on a set of objects, NLMs sequentially apply horn rules to draw further conclusions such as a property of or a relationship between objects. For example, in Blocksworld, based on a premise such as (on a b) for blocks a, b, NLMs can infer (not (clear b)). NLM has two unique features: (1) The ability to combine the predicates of different arities, and (2) size invariance & permutation equivariance on input objects, which is achieved by enumerating the permutations of input arguments.

NLM is designed to work on hand-crafted binary tensors representing propositional groundings of FOL statements. The format of these binary tensors are the multi-arity representation $(z/1, \dots, z/N)$ identical to the latent space of SimpleFOSAE.

NLM is designed for a subset of First Order Logic where every statement is a horn rule, contains no function terms (such as a function that returns an object), and all rules are applied between neighboring arities. With these assumptions, the statements fall in one of the three types below:

$$\begin{aligned} (\text{expand}) \quad & \forall x_{\# \bar{p}}; \bar{p}(X; x_{\# \bar{p}}) \leftarrow p(X), \\ (\text{reduce}) \quad & \underline{p}(X) \leftarrow \exists x_{\# p}; p(X; x_{\# p}), \\ (\text{compose}) \quad & q(X) \leftarrow \mathcal{F} \left(\bigcup_{\pi} (P \cup \underline{P} \cup \bar{P}) / \#q(\pi(X)) \right). \end{aligned}$$

Here, $p, \underline{p}, \bar{p}, q \in P, \underline{P}, \bar{P}, Q$ (respectively) are predicates, $X = (x_1, \dots)$ is a sequence of parameters, $\mathcal{F}(T)$ is a formula consisting of $\{\wedge, \vee, \neg, \top, \perp\}$ and T , $(P \cup \underline{P} \cup \bar{P}) / \#q$ is a set of predicates of the same arity as q , and $\pi(X)$ is a permutation of X , which is used for composing the predicates with the different argument orders.

All three rules can be implemented as tensor operations (Fig. 3). Given a binary tensor z/n of shape $O.^n.O \times P/n$, *expand* copies the n -th axis to $n + 1$ -th axis resulting in a shape $O.^n \top .!O \times P/n$, and *reduce* takes the max of n -th axis resulting in a shape $O.^n \top .!O \times P/n$. While the original paper proposed both min and max versions to represent \forall and \exists , with enough number of layers only one of them is necessary because $\forall x; p(\cdot, x) = \neg \exists x; \neg p(\cdot, x)$.

Finally the COMPOSE combines the information between the neighboring tensors $z/n, z/n-1, z/n+1$. In order to use the information in the neighboring arities (P, \underline{P} and \bar{P}), the input concatenates z/n with $\text{EXPAND}(z/n-1)$ and $\text{REDUCE}(z/n+1) (\rightarrow O.^n.O \times C$ where $C = P/n + P/n-1 + P/n+1)$. Next, a PERMUTE function enumerates and concatenates the results of swapping the first n axes in

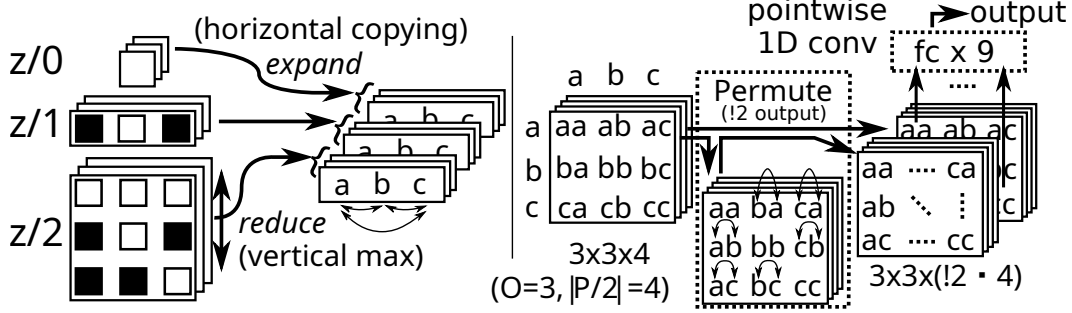


Figure 3: NLM uses EXPAND, REDUCE, PERMUTE tensor operations to combine predicates of different arities while maintaining the object-ordering / size invariance. (For matrices, PERMUTE is equivalent to transposition.)

the tensor ($\rightarrow O.^n.O \times (!n \cdot C)$). It then applies a n -D pointwise convolutional filter f_n with Q output features ($\rightarrow O.^n.O \times Q$). In the actual implementation, this n -D pointwise filter is merely a 1D convolution performed on a matrix reshaped into $O^n \times (!n \cdot C)$. It is activated by any nonlinearity σ to obtain the final result, which we denote as $\text{COMPOSE}(z, n, Q, \sigma)$. Formally, $\forall j \in 1..n, \forall o_j \in 1..O$,

$$\text{COMPOSE}(z, n, Q, \sigma)_{o_1 \dots o_n} = \sigma(f_n(\Pi_{o_1 \dots o_n})) \in \mathbb{R}^Q$$

$$\text{where } \Pi = \text{PERMUTE}\left(\text{EXPAND}(z/n-1); z/n; \text{REDUCE}(z/n+1)\right) \in \mathbb{B}^{O.^n.O \times (!n \cdot C)}$$

An NLM contains N (the maximum arity) COMPOSE operation for the neighboring arities, with appropriately omitting both ends (0 and $N + 1$) from the concatenation. We denote the result as $\text{NLM}_{Q, \sigma}(z) = (\text{COMPOSE}(z, 1, Q, \sigma), \dots, \text{COMPOSE}(z, N, Q, \sigma))$. This horizontal arity-wise compositions can be layered vertically, allowing the composition of predicates whose arities differ more than 1 (e.g., 2 layers of NLM can combine unary and quaternary predicates).

Two minor modification was made from the original paper. First, we use a slightly different tensor shapes: For the notational convenience, we use hypercube-dimensional tensors of shape $\mathbb{B}^{O.^n.O \times P/n}$, instead of the original formulation $\mathbb{B}^{O \times O-1 \times \dots \times O-n-1 \times P/n}$ which tries to reduce the size by disallowing the duplicates in the parameters. Our modification does not significantly affect the complexity of the representation because the original representation also has $O(O^n)$ complexity.

Second, we do not use the nullary predicates $z/0$ in order to disallow VAEs from encoding environment-specific information in it.

C Learning Discrete Latent STRIPS Dynamics using Back-To-Logit

$\text{APPLY}(a^i, z^{i,0})$ in AAE/Space AE is an arbitrary MLP, i.e., a neural black-box function that does not directly translates to a STRIPS action model, preventing efficient search with state-of-the-art classical planner. Cube-Space AE [Asai and Muise, 2020] addresses this issue by Back-To-Logit technique which replaces the MLP. Back-to-Logit places a so-called *cube-like graph prior* on the binary latent space / transitions. To understand the prior, the background of *cube-like graph* is necessary. This section proceeds as follows: We first explain cube-like graph (which provides the intuitive motivation) in Sec. C.1, then show the equivalence between directed cube-like graph and STRIPS (Sec. C.2). Sec. C.3 shows the equivalence between directed cube-like graph and Back-to-Logit. Finally, Sec. C.4 shows the equivalence between Back-to-Logit and discretized Kalman filter [Kalman [1960], providing an interesting connection between the classical control and STRIPS.

C.1 Cube-Like Graph

Asai and Muise [2020] identified that state transition graphs of STRIPS planning problems form a graph class called directed *cube-like graph* [Payan, 1992]. A *cube-like graph* $G(S, D) = (V, E)$ is a

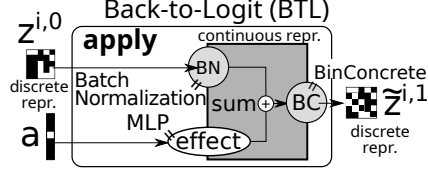


Figure 4: Back-To-Logit architecture.

simple¹ undirected graph defined by the sets S and D . Each node $v \in V$ is a finite subset of S , i.e., $v \subseteq S$. The set D is a family of subsets of S , and for every edge $e = (v, w) \in E$, the symmetric difference $d = v \oplus w = (v \setminus w) \cup (w \setminus v)$ must belong to D . For example, a unit cube is a cube-like graph because $S = \{x, y, z\}$, $V = \{\emptyset, \{x\}, \dots, \{x, y, z\}\}$, $E = \{(\emptyset, \{x\}), \dots, (\{y, z\}, \{x, y, z\})\}$, $D = \{\{x\}, \{y\}, \{z\}\}$. The set-based representation can be alternatively represented as a bit-vector, e.g., $V' = \{(0, 0, 0), (0, 0, 1), \dots, (1, 1, 1)\}$.

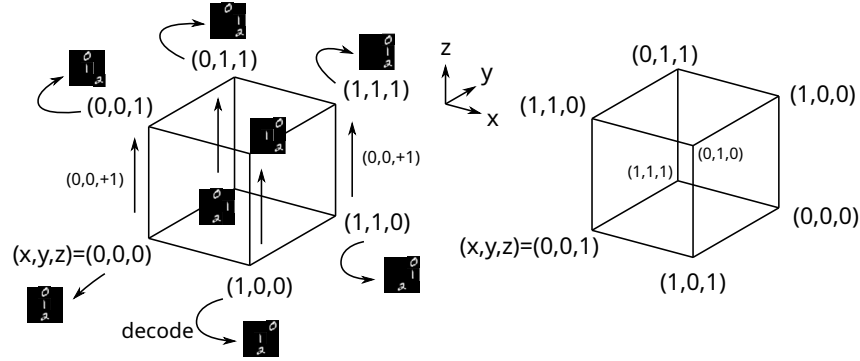


Figure 5: (Left) A graph representing a 3-dimensional cube which is a cube-like graph. (Right) A graph whose shape is identical to the left, but whose unique node embeddings are randomly shuffled.

Consider coloring a graph which forms a unit cube (Fig. 5) and whose node embeddings correspond to the latent space of some images. A cube-like graph on the left can be efficiently (i.e., by fewer colors) colored by the difference between the neighboring embeddings. Edges can be categorized into 3 labels (6 labels if directed), where each label is assigned to 4 edges which share the node embedding differences, as depicted by the upward arrows with the common node difference $(0, 0, +1)$ in the figure. This node embedding differences correspond to the set D , and each element of D represents an action, e.g., moving toward positive direction of x -axis will result in the “0” tile in the image moving to the right, and toward the positive direction of z -axis will result in the “2” tile in the image moving to the right — effectively making the actions compositional. In contrast, the graph on the right has the node embeddings that are randomly shuffled. Despite having the same topology and the same embedding size, this graph lacks the common patterns in the embedding differences like we saw on the left, thus cannot be efficiently colored by the node differences.

As such, cube-like graphs can be characterized by the smaller edge chromatic number (minimum edge coloring). Theoretical results on the graph coloring of Cube-Like Graphs guarantee the compactness of the action model.

C.2 Equivalence of Directed Cube-Like Graph and STRIPS

In STRIPS modeling, Asai and Muise [2020] used a directed version of this graph class. For every edge $e = (v, w) \in E$, there is a pair of sets $d = (d^+, d^-) = (w \setminus v, v \setminus w) \in D$ which satisfies the asymmetric difference $w = (v \setminus d^-) \cup d^+$. It is immediately obvious that this graph class corresponds to the relationship between binary states and action effects in STRIPS, $s' = (s \setminus \text{DEL}(a)) \cup \text{ADD}(a)$.

¹No duplicated edges between the same pair of nodes

It is also worth noting that in an undirected STRIPS planning problem where all actions are reversible (i.e., for any action a such that $t = a(s)$, there is an action a^{-1} that is able to reverse the effect $s = a^{-1}(t)$), the state space graph is equivalent to an undirected cube-like graph.

C.3 Back-to-Logit and its Equivalence to STRIPS

Cube-Space AE restricts the binary latent encoding and the transitions to directed cube-like graph, thereby guaranteeing the direct translation of latent space into STRIPS action model. However, since discrete representation learning is already known to be a challenge, adding a prior to it makes the training particularly difficult. Back-to-Logit (Fig. 4) was proposed in order to avoid directly operating on the discrete vectors. Instead, it converts a discrete current state $\mathbf{z}^{i,0}$ back to a continuous logit using Batch Normalization [Ioffe and Szegedy, 2015, BN], takes a sum with a continuous *effect* produced by an additional MLP $\text{EFFECT}(\mathbf{a}^i)$, and re-discretize the resulting logit using Binary Concrete.

$$\tilde{\mathbf{z}}^{i,1} = \text{APPLY}(a, \mathbf{z}^{i,0}) = \text{BC}(\text{BN}(\mathbf{z}^{i,0}) + \text{EFFECT}(\mathbf{a}^i)).$$

Implementation note: Since \mathbf{a}^i is a probability vector over A action ids, the additional MLP can be merely an embedding, i.e. $\text{EFFECT}(\mathbf{a}^i) = \mathbf{E}\mathbf{a}^i$, where \mathbf{E} is an embedding matrix. It also helps the training if we apply BatchNorm on the effect vector. Therefore, a possible implementation is

$$\text{APPLY}(a, \mathbf{z}^{i,0}) = \text{BC}(\text{BN}(\mathbf{z}^{i,0}) + \text{BN}(\mathbf{E}\mathbf{a}^i)).$$

States learned by BTL has the following property:

Theorem 1. [Asai and Muise, 2020] *Under the same action a , state transitions are bitwise monotonic, deterministic, and restricted to three mutually exclusive modes, i.e., for each bit j :*

$$\begin{aligned} (\text{add:}) \quad & \forall i; (\mathbf{z}_j^{i,0}, \mathbf{z}_j^{i,1}) \in \{(0, 1), (1, 1)\} \text{ i.e. } \mathbf{z}_j^{i,0} \leq \mathbf{z}_j^{i,1} \\ (\text{del:}) \quad & \forall i; (\mathbf{z}_j^{i,0}, \mathbf{z}_j^{i,1}) \in \{(1, 0), (0, 0)\} \text{ i.e. } \mathbf{z}_j^{i,0} \geq \mathbf{z}_j^{i,1} \\ (\text{nop:}) \quad & \forall i; (\mathbf{z}_j^{i,0}, \mathbf{z}_j^{i,1}) \in \{(0, 0), (1, 1)\} \text{ i.e. } \mathbf{z}_j^{i,0} = \mathbf{z}_j^{i,1} \end{aligned}$$

This theorem guarantees that each action deterministically sets a certain bit on and off in the binary latent space. Therefore, the actions and the transitions satisfy the STRIPS state transition rule $s' = (s \setminus \text{DEL}(a)) \cup \text{ADD}(a)$, thus enabling a direct translation from neural network weights to PDDL modeling language.

The proof is straightforward from the monotonicity of the BatchNorm and Binary Concrete. Note that we assume BatchNorm’s additional scale parameter γ is kept positive or disabled.

Proof. For readability, we omit j and assumes a 1-dimensional case. Let $e = \text{EFFECT}\mathbf{a} \in \mathbb{R}$. Note that e is a constant for the fixed input \mathbf{a} . At the limit of annealing, Binary Concrete BC becomes a STEP function, which is also monotonic. BN is monotonic because we assumed the scale parameter γ of BN is positive, and the main feature of BN also only scales the variance of the batch, which is always positive. Then we have

$$\mathbf{z}^{i,1} = \text{STEP}(\text{BN}(\mathbf{z}^{i,0}) + e).$$

The possible values a pair $(\mathbf{z}^{i,0}, \mathbf{z}^{i,1})$ can have is $(0, 0), (0, 1), (1, 0), (1, 1)$. Since both STEP and BN are deterministic at the testing time (See Ioffe and Szegedy [2015]), we consider the deterministic mapping from $\mathbf{z}^{i,0}$ to $\mathbf{z}^{i,1}$. There are only 4 deterministic mappings: $\{(0, 1), (1, 1)\}$, $\{(1, 0), (0, 0)\}$, $\{(0, 0), (1, 1)\}$, and lastly $\{(0, 1), (1, 0)\}$. Thus our goal is now to show that the last mapping is impossible in latent space $\{\dots(\mathbf{z}^{i,0}, \mathbf{z}^{i,1})\dots\}$.

To prove this, first, assume $(\mathbf{z}^{i,0}, \mathbf{z}^{i,1}) = (0, 1)$ for some index i . Then

$$1 = \text{STEP}(\text{BN}(0) + e). \Rightarrow \text{BN}(0) + e > 0. \Rightarrow \text{BN}(1) + e > 0. \Rightarrow \forall i; \text{BN}(\mathbf{z}^{i,0}) + e > 0.$$

The second step is due to the monotonicity $\text{BN}(0) < \text{BN}(1)$. This shows $\mathbf{z}^{i,1}$ is constantly 1 regardless of $\mathbf{z}^{i,0}$, therefore it proves that $(\mathbf{z}^{i,0}, \mathbf{z}^{i,1}) = (1, 0)$ cannot happen in any i .

Likewise, if $(\mathbf{z}^{i,0}, \mathbf{z}^{i,1}) = (1, 0)$ for some index i ,

$$0 = \text{STEP}(\text{BN}(1) + e). \Rightarrow \text{BN}(1) + e < 0. \Rightarrow \text{BN}(0) + e < 0. \Rightarrow \forall i; \text{BN}(\mathbf{z}^{i,0}) + e < 0.$$

Therefore, $z^{i,1} = 0$ regardless of $z^{i,0}$, and thus $(z^{i,0}, z^{i,1}) = (0, 1)$ cannot happen in any i .

Finally, if the data points do not contain $(0, 1)$ or $(1, 0)$, then by assumption they do not coexist. Therefore, the embedding learned by BTL cannot contain $(0, 1)$ and $(1, 0)$ at the same time. \square

C.4 BTL is Equivalent to a Discretized Kalman Filter

Kalman filter Kalman [1960] is a linear dynamics model for noisy observation. Given a current state vector \mathbf{x}_t , a control vector \mathbf{u}_t , and a gaussian noise vector \mathbf{w}_t , the continuous state dynamics is modeled as

$$\mathbf{x}_{t+1} = \mathbf{F}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \mathbf{w}_t,$$

assuming a static state dynamics \mathbf{F} and a static control-input model \mathbf{B} .

We notice that our BTL formulation is equivalent to a discretization of Kalman filter. This equivalence establishes an interesting connection between STRIPS and classical control. Compare Kalman Filter with our BTL model below, where we additionally expanded Binary-Concrete BC to a sigmoid σ , annealing parameter τ , and a sample \mathbf{l} from LOGISTIC(0,1) distribution:

$$\tilde{z}^{i,1} = \text{BC}(\text{BN}(\mathbf{z}^{i,0}) + \text{EFFECT}(\mathbf{a}^i)) = \sigma((\text{BN}(\mathbf{z}^{i,0}) + \text{EFFECT}(\mathbf{a}^i) + \mathbf{l})/\tau).$$

Since Batch Normalization BN only applies a shift β and a scale γ to each element of the vector $\mathbf{z}^{i,0}$, it is a linear operation that can be represented by a matrix \mathbf{W} . We also noted that EFFECT can be an arbitrary MLP, thus can be an embedding matrix \mathbf{E} and BatchNorm. Again, since BatchNorm is equivalent to a matrix, we can represent the effect embedding as

$$\text{EFFECT}(\mathbf{a}^i) = \text{BN}(\mathbf{E}\mathbf{a}^i) = \mathbf{E}'\mathbf{a}^i.$$

As such, we can reformulate BTL as

$$\tilde{z}^{i,1} = \sigma((\mathbf{W}\mathbf{z}^{i,0} + \mathbf{E}'\mathbf{a}^i + \mathbf{l})/\tau).$$

The limit of this model at $\tau \rightarrow 0$ is a step function, thus our model is a discretization of Kalman filter with a logistic noise, i.e.,

$$\tilde{z}^{i,1} = \text{STEP}(\mathbf{W}\mathbf{z}^{i,0} + \mathbf{E}'\mathbf{a}^i + \mathbf{l}).$$

C.5 Effect Rule Extraction

To extract the effects of an action \mathbf{a} from Cube-Space AE, we compute $\text{ADD}(\mathbf{a}) = \text{APPLY}(\mathbf{a}, \mathbf{0})$ and $\text{DEL}(\mathbf{a}) = 1 - \text{APPLY}(\mathbf{a}, \mathbf{1})$ for each action \mathbf{a} , where $\mathbf{0}, \mathbf{1}$ are vectors filled by zeros/ones and has the same size as the binary embedding. Since APPLY deterministically sets values to 0 or 1, feeding these vectors is sufficient to see which bit it turns on and off. For each j -th bit that is 1 in each result, a corresponding proposition is added to the add/delete-effect, respectively.

In FOSAE++, we extract the effects from parameter-bounded subspace $\mathbf{z}_{\dagger}^{i,0}, \mathbf{z}_{\dagger}^{i,1}$. The representation is a tuple $\mathbf{z}_{\dagger} = (z_{\dagger}/1 \cdots z_{\dagger}/N)$, where $z_{\dagger}/n \in \mathbb{B}^{\#a \times ?n \times \#a \times P/n}$. BTL then operates on flattened and concatenated binary vectors of size $\sum_n \#a^n P/n$: The input, the output, and the effect share this shape. We extract the effects from this BTL vector in the same manner as noted above. After the extraction, however, each bit is converted to a lifted predicate according to the position. For example, when a bit that corresponds to $z_{\dagger}/2_{1,2,5}$ has turned from 0 to 1, then the add-effect contains $\text{p5}(\text{?arg1}, \text{?arg2})$, where ?arg1 is a parameter used in the lifted PDDL encoding.

We show an example of such a learned PDDL model below, obtained from 8-Puzzle with $P/1 = P/2 = 333$ (reformatted for readability). Note that we disabled the nullary predicates $P/0$ and $z/0$, which consumes the first 333 dimensions in the flattened vector. Another note is that we also count the number of appearances of each action in the training dataset. If an action label is never used in the dataset, it is not exported in the resulting PDDL output. Thus, the index for the action starts from 7 in the example.

```
(define (domain latent)
  (:requirements :strips :negative-preconditions)
  (:predicates
    (p333 ?x0) ... (p665 ?x0)
    (p666 ?x0 ?x1) ... (p998 ?x0 ?x1))
```

```

(:action a7 :parameters (?x0) :precondition
  (and (p339 ?x0) (p388 ?x0) (p391 ?x0) (p398 ?x0) (p402 ?x0)
        (p420 ?x0) (p421 ?x0) (p446 ?x0) (p447 ?x0) (p473 ?x0)
        (p475 ?x0) (p489 ?x0) (p491 ?x0) (p502 ?x0) (p516 ?x0)
        (p559 ?x0) (p588 ?x0) (p615 ?x0) (p641 ?x0) (p648 ?x0)
        (p831 ?x0 ?x0) (p950 ?x0 ?x0) (not (p333 ?x0))
        (not (p371 ?x0)) (not (p375 ?x0)) (not (p388 ?x0))
        (not (p402 ?x0)) (not (p406 ?x0)) (not (p421 ?x0))
        (not (p447 ?x0)) (not (p454 ?x0)) (not (p504 ?x0))
        (not (p508 ?x0)) (not (p519 ?x0)) (not (p524 ?x0))
        (not (p526 ?x0)) (not (p562 ?x0)) (not (p584 ?x0))
        (not (p593 ?x0)) (not (p617 ?x0)) (not (p640 ?x0))
        (not (p652 ?x0)) (not (p824 ?x0 ?x0)) (not (p892 ?x0 ?x0))
        (not (p926 ?x0 ?x0)) (not (p975 ?x0 ?x0))
        (not (p994 ?x0 ?x0)))
  :effect
  (and (p349 ?x0) (p361 ?x0) (p366 ?x0) (p370 ?x0) (p371 ?x0)
        (p378 ?x0) (p381 ?x0) (p385 ?x0) (p388 ?x0) (p401 ?x0)
        (p408 ?x0) (p421 ?x0) (p432 ?x0) (p445 ?x0) (p454 ?x0)
        (p475 ?x0) (p491 ?x0) (p496 ?x0) (p502 ?x0) (p503 ?x0)
        (p504 ?x0) (p507 ?x0) (p517 ?x0) (p526 ?x0) (p550 ?x0)
        (p562 ?x0) (p563 ?x0) (p575 ?x0) (p584 ?x0) (p588 ?x0)
        (p599 ?x0) (p601 ?x0) (p607 ?x0) (p612 ?x0) (p617 ?x0)
        (p631 ?x0) (p640 ?x0) (p641 ?x0) (p647 ?x0) (p656 ?x0)
        (p663 ?x0) (p724 ?x0 ?x0) (p768 ?x0 ?x0) (p831 ?x0 ?x0)
        (p902 ?x0 ?x0) (p911 ?x0 ?x0) (p993 ?x0 ?x0) (not (p339 ?x0))
        (not (p355 ?x0)) (not (p365 ?x0)) (not (p391 ?x0))
        (not (p397 ?x0)) (not (p398 ?x0)) (not (p402 ?x0))
        (not (p406 ?x0)) (not (p422 ?x0)) (not (p446 ?x0))
        (not (p447 ?x0)) (not (p448 ?x0)) (not (p451 ?x0))
        (not (p456 ?x0)) (not (p472 ?x0)) (not (p473 ?x0))
        (not (p478 ?x0)) (not (p489 ?x0)) (not (p490 ?x0))
        (not (p495 ?x0)) (not (p516 ?x0)) (not (p518 ?x0))
        (not (p524 ?x0)) (not (p525 ?x0)) (not (p527 ?x0))
        (not (p534 ?x0)) (not (p544 ?x0)) (not (p559 ?x0))
        (not (p561 ?x0)) (not (p615 ?x0)) (not (p624 ?x0))
        (not (p629 ?x0)) (not (p642 ?x0)) (not (p646 ?x0))
        (not (p651 ?x0)) (not (p653 ?x0)) (not (p720 ?x0 ?x0))
        (not (p813 ?x0 ?x0)) (not (p824 ?x0 ?x0)) (not (p892 ?x0 ?x0))
        (not (p894 ?x0 ?x0)) (not (p926 ?x0 ?x0)) (not (p931 ?x0 ?x0))
        (not (p975 ?x0 ?x0)) (not (p994 ?x0 ?x0)))
  (:action a22 :parameters (?x0) :precondition
    ...

```

C.6 Precondition Learning with Dynamics Reversed in Time

In the main text, we simplified the model by showing only the forward dynamics, i.e., the dynamics in the same direction as the time. This forward dynamics can model the effects (add/delete) of the actions. However, the forward dynamics is insufficient for learning the preconditions of the actions. The original CSAE paper [Asai and Muise, 2020] used an add-hoc method that extracts common bits of the current states.

In contrast, we added a network that uses the same BTL mechanism that is applied backward in time, i.e., predict the current state $z^{i,0}$ from a successor state $z^{i,1}$ and a one-hot action vector a^i . We named the network $\text{REGRESS}(z^{i,1}, a^i)$, alluding to the *regression planning* [Alcázar et al., 2013] literature.

In REGRESS, add-effects and delete-effects now correspond to *positive preconditions* and *negative preconditions*. A positive precondition (normal precondition) requires that a proposition is \top prior to using an action. In contrast, a negative precondition requires that a proposition is \perp prior to using an action. While negative preconditions are (strictly speaking) out of STRIPS formalism, it is commonly supported by the modern classical planners that participated in the recent competitions. To extract the preconditions from the network, we can use the same method used for extracting the effects from the progressive/forward dynamics.

The entire model thus looks as follows.

(encoder)	$z^{i,0}, z^{i,1} = \text{ENCODE}(o^{i,0}), \text{ENCODE}(o^{i,1})$
(action parameters, uses NLMs)	$x^i = \text{PARAMS}(z^{i,0}, z^{i,1})$
(parameter-bound subspace extraction)	$z_{\dagger}^{i,0}, z_{\dagger}^{i,1} = \text{BIND}(z^{i,0}, x^i), \text{BIND}(z^{i,1}, x^i)$
(action assignment)	$a^i = \text{ACTION}(z_{\dagger}^{i,0}, z_{\dagger}^{i,1})$
(bounded forward dynamics)	$\tilde{z}_{\dagger}^{i,1} = \text{APPLY}(z_{\dagger}^{i,0}, a^i)$
(bounded backward dynamics)	$\tilde{z}_{\dagger}^{i,0} = \text{REGRESS}(z_{\dagger}^{i,1}, a^i)$
(reflection to global forward dynamics)	$\tilde{z}^{i,1} = z^{i,0} - \text{UNBIND}(z_{\dagger}^{i,0}, x^i) + \text{UNBIND}(\tilde{z}_{\dagger}^{i,1}, x^i)$
(reflection to global backward dynamics)	$\tilde{z}^{i,0} = z^{i,1} - \text{UNBIND}(z_{\dagger}^{i,1}, x^i) + \text{UNBIND}(\tilde{z}_{\dagger}^{i,0}, x^i)$
(reconstructions)	$\tilde{o}^{i,0}, \tilde{o}^{i,1} = \text{DECODE}(z^{i,0}), \text{DECODE}(z^{i,1})$
(reconstruction based on forward dynamics)	$\tilde{o}^{i,1} = \text{DECODE}(\tilde{z}^{i,1})$
(reconstruction based on backward dynamics)	$\tilde{o}^{i,0} = \text{DECODE}(\tilde{z}^{i,0})$

The total loss is $\ell(o^{i,0}, \tilde{o}^{i,0}) + \ell(o^{i,1}, \tilde{o}^{i,1}) + \ell(o^{i,1}, \tilde{o}^{i,1}) + \ell(o^{i,0}, \tilde{o}^{i,0}) + \ell(z^{i,1}, \tilde{z}^{i,1}) + \ell(z_{\dagger}^{i,1}, \tilde{z}_{\dagger}^{i,1}) + \ell(z^{i,0}, \tilde{z}^{i,0}) + \ell(z_{\dagger}^{i,0}, \tilde{z}_{\dagger}^{i,0}) + \text{Reg}$.

D Implementation Detail

We based our code on a publicly-available Latplan source code repository (<https://github.com/guicho271828/latplan/>), which is based on Keras Deep Learning library. The repository hosts its own Genetic Algorithm based hyperparameter tuner which we mention several times later.

D.1 Input Data Format and Loss Functions

As we discussed in the earlier appendix section, our final loss function consists of 9 terms: $\ell(\mathbf{o}^{i,0}, \tilde{\mathbf{o}}^{i,0}) + \ell(\mathbf{o}^{i,1}, \tilde{\mathbf{o}}^{i,1}) + \ell(\mathbf{o}^{i,1}, \tilde{\mathbf{o}}^{i,1}) + \ell(\mathbf{o}^{i,0}, \tilde{\mathbf{o}}^{i,0}) + \ell(\mathbf{z}^{i,1}, \tilde{\mathbf{z}}^{i,1}) + \ell(\mathbf{z}_{\dagger}^{i,1}, \tilde{\mathbf{z}}_{\dagger}^{i,1}) + \ell(\mathbf{z}^{i,0}, \tilde{\mathbf{z}}^{i,0}) + \ell(\mathbf{z}_{\dagger}^{i,0}, \tilde{\mathbf{z}}_{\dagger}^{i,0}) + \text{Reg}$. We first describe the input data format that is shared among different domains, then the loss functions defined on it.

Our input/output format $\mathbf{o} \in \mathbb{R}^{O \times F}$ consists of O objects (environment-dependent) each having F features. F features consists of image-based features and the coordinate / dimension-based features (Fig. 6). All image patches extracted from the observations are resized into a fixed height H , width W and color channel $C = 3$. Each flattened object vector has the size $F = H \times W \times C + 4$. The last 4 dimensions contain the center coordinate and the actual height/width before the resizing. Out of the 9 terms in the total loss, the terms that apply to the object vectors of this form are $\ell(\mathbf{o}^{i,0}, \tilde{\mathbf{o}}^{i,0})$, $\ell(\mathbf{o}^{i,1}, \tilde{\mathbf{o}}^{i,1})$, $\ell(\mathbf{o}^{i,1}, \tilde{\mathbf{o}}^{i,1})$, $\ell(\mathbf{o}^{i,0}, \tilde{\mathbf{o}}^{i,0})$.



Figure 6: Data array representing a single object.

For this data format, the loss function consists of the mean square value of the image part, and the square sum of the coordinate / dimension parts. We do not average the losses for the coordinates/dimensions to avoid making the gradient minuscule. To further enhance this direction, we additionally have the *coordinate loss amplifier* λ tuned by GA, as it is often the case that the object location has more visual impact on the reconstruction. Note that, for the tuning with the validation set and the evaluation with the test set, we set $\lambda = 1$ in order to obtain the consistent, comparable measuring. λ is altered only during the training. Formally, for the i -th objects in the input \mathbf{o} and the reconstruction $\tilde{\mathbf{o}}$, we define the loss as follows. These losses are averaged over the objects and the batch dimensions.

$$\ell(\mathbf{o}_i, \tilde{\mathbf{o}}_i) = \frac{1}{HWC} \|\mathbf{o}_{i,1..HWC} - \tilde{\mathbf{o}}_{i,1..HWC}\|_2^2 + \lambda \|\mathbf{o}_{i,HWC..F} - \tilde{\mathbf{o}}_{i,HWC..F}\|_2^2$$

We call the remaining losses except the regularization terms as “latent dynamics loss”. They operate on the binary latent data activated by BinaryConcrete: $\ell(\mathbf{z}^{i,1}, \tilde{\mathbf{z}}^{i,1})$, $\ell(\mathbf{z}_{\dagger}^{i,1}, \tilde{\mathbf{z}}_{\dagger}^{i,1})$, $\ell(\mathbf{z}^{i,0}, \tilde{\mathbf{z}}^{i,0})$, $\ell(\mathbf{z}_{\dagger}^{i,0}, \tilde{\mathbf{z}}_{\dagger}^{i,0})$. However, note that during the training, all values used for the loss calculation are still continuous due to BinaryConcrete’s annealing. This means we can’t use Binary Cross Entropy BCE, the standard loss function for binary classifications, because the “training data” is also a noisy probability value. It is also in fact symmetric — as we discussed in the previous appendix sections, the role of these losses is not only just to obtain the accurate dynamics, but *also to shape the state representation toward a cube-like graph*. While there are several candidate loss functions that can be considered, we adapted Symmetric Cross Entropy [Wang et al., 2019] designed for noisy labels, which simply applies BCE in both ways. Formally, given $\text{BCE}(\mathbf{z}, \tilde{\mathbf{z}}) = -\sum_i (\mathbf{z}_i \log \tilde{\mathbf{z}}_i + (1 - \mathbf{z}_i) \log(1 - \tilde{\mathbf{z}}_i))$,

$$\ell(\mathbf{z}, \tilde{\mathbf{z}}) = \text{BCE}(\mathbf{z}, \tilde{\mathbf{z}}) + \text{BCE}(\tilde{\mathbf{z}}, \mathbf{z}).$$

The remaining regularization losses include KL divergence for Discrete VAEs, as well as the L1 regularization for the latent vectors $\mathbf{z}^{i,0}$, $\mathbf{z}^{i,1}$ which were proven useful in Asai and Kajino [2019].

Finally, we define the magnitude and the warmup of each loss. The magnitude multiplies each loss and is tuned by the GA tuner. Similar to λ , the values are set to 1 during the evaluation. We have α for the L1 regularization, β for KL divergence, and γ for the latent dynamics loss.

The warmup mechanism works by setting these values to 0 until the training reaches a certain epoch defined by the ratio r relative to the total number of epochs. We used the warmup r_α , r_γ for α , γ , as well as r_{rec} and the r_{dyn} for the main reconstruction loss $\ell(\mathbf{o}^{i,0}, \tilde{\mathbf{o}}^{i,0}) + \ell(\mathbf{o}^{i,1}, \tilde{\mathbf{o}}^{i,1})$ and the dynamics-based reconstruction loss $\ell(\mathbf{o}^{i,1}, \tilde{\mathbf{o}}^{i,1}) + \ell(\mathbf{o}^{i,0}, \tilde{\mathbf{o}}^{i,0})$.

The motivation behind these magnitudes and warmups is to balance the speed of convergence of the various parts of the networks. Depending on the hyperparameters (depth, width of the layers), occasionally the network completely ignores the dynamics, falling into something similar to (but the mechanism will be very different from) a mode collapse where the effect caused by the dynamics is empty, e.g., the forward dynamics produces the same state as the current state.

Another failure mode of the network is that the dynamics loss is too strong, due to its BCE which could become too large compared to the reconstruction loss. As a result, the network learns a perfect but meaningless latent space dynamics which does not produce correct reconstructions. Tuning the warmup and balancing the losses addressed these issues.

D.2 Network Detail

FOSAE++ consists of 6 networks: ENCODE, PARAMS, ACTION, APPLY, REGRESS, DECODE. Note that BIND and UNBIND are weight-less operations.

D.2.1 Encoder

ENCODE can be divided into trivial continuous feature extraction phase (pre-encoder) and the actual FOSAE encoder. The feature extractor is a simple 1D pointwise convolution over the objects, i.e., same as applying the same Dense/Fully-Connected(FC) layer on each individual object. Its depth and the hidden layer width is tuned by the GA tuner. All activations except the last Binary Concrete are Rectified Linear Unit [Fukushima, 1980, Nair and Hinton, 2010]. The architecture is illustrated in Fig. 6.

We should note that we assign the same number of predicates to different arities: $P/1 = P/2 \dots = P/N$. In this Network Detail section, we denote $\bar{P} = P/1 \dots = P/N$ as a hyperparameter that specifies the number, although P was already used in the main text as the total number of predicates (i.e., $\sum_n P/n$).

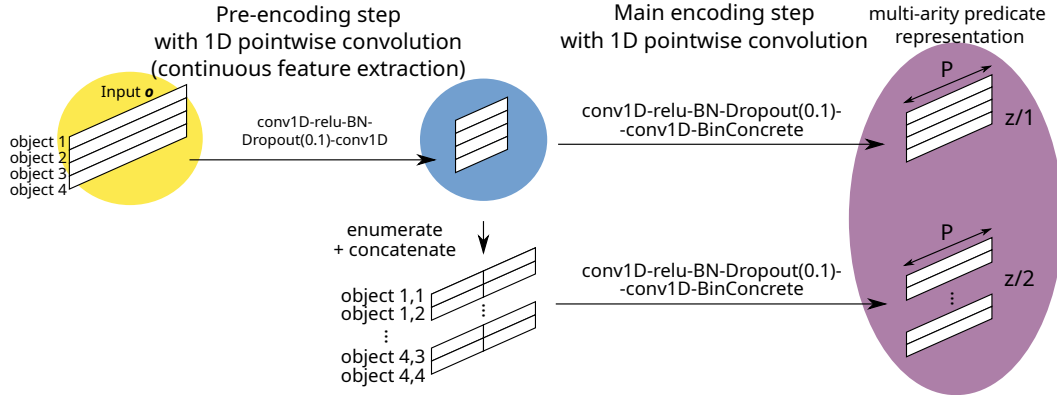


Figure 7: The encoder has two stages. The first stage extracts a compact but still continuous representation of each object. From this compressed representation, the second stage identifies properties of objects as well as relationships/predicates between objects.

D.2.2 Params

PARAMS consists of multiple NLM layers, as depicted in Fig. 8. All hidden activations are ReLU, and the width Q of PARAMS is tuned by the GA tuner.

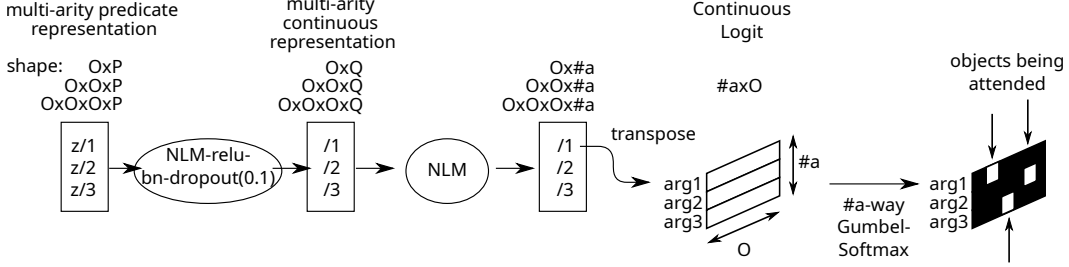


Figure 8: PARAMS network using NLMs.

D.2.3 Action, Apply / Regress

Since ACTION takes parameter-bounded representations $(z_{\dagger}^{i,0}, z_{\dagger}^{i,1})$ whose objects are already selected and appropriately reordered by BIND, we simply apply an MLP whose last layer has an output size A , and is activated by Gumbel-Softmax. This results in selecting one action, represented by a one-hot vector. Before applying the MLP, we should flatten and concatenate the tuple of vectors, each of size $\#a \times \text{?} \times \#a \times P$. The depth and the hidden layer width of the MLP is automatically tuned.

For APPLY and REGRESS, we similarly flatten the input and directly apply the BTL structure described in the earlier sections. The EFFECT network in $\tilde{z}^{i,1} = \text{BC}(\text{BN}(z^{i,0}) + \text{EFFECT}(a^i))$ is a single linear layer without bias combined with a batch normalization. In other words, with the embedding matrix E ,

$$\tilde{z}^{i,1} = \text{BC}(\text{BN}(z^{i,0}) + \text{BN}(Ea^i)).$$

D.2.4 Decoder

The decoder consists of NLMs followed by a post-decoder that shares the same width and depth as the pre-encoder. In Blocksworld, we additionally made the decoder a Bayesian Neural Network to absorb the uncertainty in the reconstruction. Details are available in (Sec. F.2).

D.3 Hyperparameter Tuner

The tuning system assumes that the hyperparameter configuration is a vector/dict of categorical/finite candidates. The tuner is a textbook integer GA with uniform crossover and point mutation. Assume that a hyperparameter configuration is represented by H values. A new configuration $p = \{p_1, \dots, p_H\}$ is created from two parents $q = \{q_1, \dots, q_H\}$ and $r = \{r_1, \dots, r_H\}$ by $\forall i; p_i = \text{RandomChoice}(q_i, r_i)$, and then a single value is randomly mutated, i.e., for $m = \text{RandomChoice}(1..H)$, $p_m \leftarrow \text{RandomChoice}(\text{ValidValuesOf}(p_m) \setminus \{p_m\})$. It stores all evaluated configurations and never re-evaluates the same configuration.

In the beginning of the tuning process, it bootstraps the initial population with a certain number of configurations. New configurations are evaluated by the validation loss, then pushed in a sorted list. A certain number of best-performing configurations in the list are considered as a “live” population. In each iteration, it selects two parents from the live population by inverse weighted random sampling of the score, preferring the smaller validation losses but also occasionally selecting a second-tier parent. Non-performing configurations will “die” by being pushed away from the top as the algorithm finds better configurations. The evaluation and insertion to the queue is asynchronous and all processes can run in parallel.

E Training Details and Hyperparameters

We trained multiple network configurations on a distributed compute cluster equipped with Tesla K80 and Xeon E5-2600 v4, which is a rather old hardware. The list of hyperparameters are shown in Table 3.

In all experiments, we used the total of 5000 state transitions (10000 states) from the training environments. The details of data collection for each domain is available in the later sections. This dataset is divided into training/validation/test sets (90%:5%:5%). The Genetic Algorithm hyperparameter tuner uses the validation loss as the evaluation metric.

We set a limit of 1500 total runs for each environment, with 100 initial population, and ran maximum 100 processes in parallel for each environment. As an additional trick, to avoid testing unpromising candidates (e.g., those with diverging loss), the epoch parameter is forced to be 50 in the first 100 configurations and the runs finish quickly. The rest of the runs use these initial populations as the parents, but replaces the epoch with an appropriate value selected from the candidates.

Training parameters	
Optimizer	Rectified Adam, RMSProp, Nadam, Adam
Epochs	100, 333, 1000
Batch size	100, 333, 1000
Learning rate	$10^{-2}, \dots, 10^{-5}$
Gradient norm clipping	0.1, 1.0, 10
Initial annealing temperature τ_{\max}	10, 5, 2
Final annealing temperature τ_{\min}	0.2, 0.5, 0.7
Coordinate loss amplifier	1, 10, 100
Network shape parameters	
#p	1, 2, 3
#a	1, 2, 3
P	10, 33, 100, 333, 1000
A	10, 33, 100, 333, 1000
Pre-encoder/decoder output dimension	10, 33, 100, 333, 1000
Encoder hidden layer width	10, 33, 100, 333, 1000
Encoder depth	2
Decoder hidden layer width	10, 33, 100, 333, 1000
Decoder depth	2
Params/Action/Apply/Regress hidden layer width	10, 33, 100, 333, 1000
Params/Action/Apply/Regress depth	1, 2, 3
Params/Action/Apply/Regress hidden activation	ReLU, tanh
Regularization Parameters	
Latent L1 regularization α	0.0, 0.01, 0.1, 0.2
KL divergence coefficient β	0.0, 0.1, 0.3
Latent dynamics loss γ	0.1, 0.2, 0.3, 0.5, 0.8
Warmup parameters	
Latent L1 regularization d_{α}	0.0, 0.01, 0.1, 0.2, 0.3, 0.5, 0.8
Main reconstruction loss d_{rec}	0.0, 0.01, 0.1, 0.2, 0.3, 0.5, 0.8
Dynamics-based reconstruction loss d_{dyn}	0.0, 0.01, 0.1, 0.2, 0.3, 0.5, 0.8
Latent dynamics loss d_{γ}	0.0, 0.01, 0.1, 0.2, 0.3, 0.5, 0.8
Early stopping	
Explosion detection	$10 \times$ the loss value at the epoch 0

Table 3: List of hyperparameters tuned by the Genetic Algorithm.

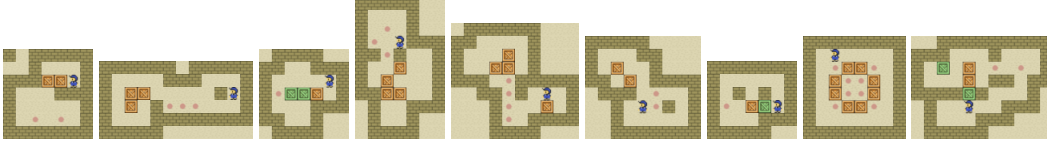


Figure 9: Visualized sokoban problems. The first 5 are used for the training, and the rest are used for evaluation. Pink dots depict the goals that the player pushes the blocks onto. Green boxes are already on one of the goals.

F Domain-Wise Details

F.1 Sokoban

We generated 10000 transitions of each training problem using Dijkstra search from the initial state. We shuffled the 50000 transitions, subsampled 5000 transitions out of 50000, then stored them in a single archive. The rendering and other data are obtained from PDDL Gym library Silver and Chitnis [2020]. Each tile is resized into 16x16 pixels, and the tile ordering is also shuffled.

In order to make the training data dimension consistent and make it convenient for GPU-based training, we performed a so-called *Random Object Masking* which removes a certain number of objects from each state via random selection. The idea is similar to masking-based image augmentation commonly used in the computer vision literature, but the purpose is different from those: Ours has more emphasis on having the consistent number of objects in each state. For example, Sokoban training problem 0 (leftmost in Fig. 9) has 49 tiles, while the problem 1 (the second picture) has 72 tiles. In the combined archive, the number of objects is set to the smallest among the problems.

We also removed certain tiles that cannot be reached by the player. For example, in problem 0, the three floors on the top left corner cannot be reached by the player. Similarly in problem 1, the bottom right corner is not reachable by the player. We performed a simple custom reachability analysis using the meta-information obtained from PDDL Gym library. This helped reducing the dataset size and thus the training.

Finally, during the dataset merging, we accounted for the potential location bias caused by the map size difference. For example, if we preserved the original x, y locations, the tiles tend to be biased around 0, 0 and the location around $(x, y) = (12, 12)$ (by tiles) is never occupied by any tile. To address the issue, for each state pair, we relocated the entire environment by a value selected uniformly randomly from a certain width. The width is decided from the maximum dimension of all training problems, i.e., 12x12. For example, a state pair in problem 4 (which has a 10x9 map dimension) will be shifted by a random value between 0..(12 - 10) in x -axis, and 0..(12 - 9) in y -axis.

F.2 Blocksworld

We generated 5000 random state transitions using Photorealistic-Blocksworld dataset Asai [2018], which in turn is based on CLEVR Johnson et al. [2017] dataset generator. It uses Blender 3D rendering engine to produce realistic effects such as metallic reflection, surface materials and shadows. The objects are extracted from the information available from the generator metadata. We cropped the image region suggested by the metadata, resized them into 32x32x3 image patches and stored them in an archive. The size of the objects reported by the metadata may vary and is noisy due to the camera jitter, object location, and the heuristics used for extracting the objects. This is a case even for the objects that are not moved by the random actions.

Each transition is generated by sampling the current state and then randomly moving a block. To sample a state, we first generate a set of block configurations (color, material, shape, size), then placing them randomly on a straight line in the 3D environment without collisions. When we move a block, we select a set of blocks of which nothing else is on top, choose one randomly, pick a new horizontal location and place it at the lowest non-colliding height. We ensure that the block is always moved, i.e., not stacked on top of the same block, and not on the floor if it is already on the floor.

In Blocksworld, we noticed that a same conceptual symbolic action (such as move) may have a non-deterministic outcome in the image space while each individual concept is discrete and deterministic. For example, move may relocate a block onto a floor, but the resulting position is chosen randomly during the dataset generation, i.e., it can be placed anywhere on the floor.

To model this uncertainty in our framework, we used a Bayesian Neural Network layers in the decoder for Blocksworld domain. The final output $\mathbf{o} \in \mathbb{R}^{O \times F}$ is produced by two NLMs of output feature size F , each producing the mean and the variance vectors $\mu, \epsilon \in \mathbb{R}^{O \times F}$, and the output reconstruction is generated by the random sampling: $\mathbf{o} = \mu + \sigma \cdot \epsilon$, where ϵ is a random noise vector following normal distribution $\mathcal{N}(0, 1)$. During the testing (e.g., visualization), the random sampling is disabled and we use the value μ for rendering.

F.3 8-Puzzle

The 8-puzzle domain generator is directly included in the Latplan code base. It uses MNIST dataset for its tiles and the tile size is set to 16. Similarly, we generated 5000 random state transitions using the generator.

G PDDL generation and Planning Experiments

G.1 Example output: 8-Puzzle

