

# Symbolic Network: Generalized Neural Policies for Relational MDPs

Sankalp Garg<sup>1</sup>, Aniket Bajpai<sup>1</sup>, Mausam<sup>1</sup>

<sup>1</sup>Indian Institute of Technology Delhi

{sankalp262199, quantum.computing96}@gmail.com, mausam@cse.iitd.ac.in

## Abstract

A Relational Markov Decision Process (RMDP) is a first-order representation to express all instances of a single probabilistic planning domain with possibly unbounded number of objects. Early work in RMDPs outputs generalized (instance-independent) first-order policies or value functions as a means to solve *all* instances of a domain at once. Unfortunately, this line of work met with limited success due to inherent limitations of the representation space used in such policies or value functions. Can neural models provide the missing link by easily representing more complex generalized policies, thus making them effective on all instances of a given domain?

We present SYMNET, the first neural approach for solving RMDPs that are expressed in the probabilistic planning language of RDDDL. SYMNET trains a set of shared parameters for an RDDDL domain using training instances from that domain. For each instance, SYMNET first converts it to an instance graph and then uses relational neural models to compute node embeddings. It then scores each ground action as a function over the first-order action symbols and node embeddings related to the action. Given a new test instance from the same domain, SYMNET architecture with pre-trained parameters scores each ground action and chooses the best action. This can be accomplished in a single forward pass *without any retraining* on the test instance, thus implicitly representing a neural generalized policy for the whole domain. Our experiments on nine RDDDL domains from IPPC demonstrate that SYMNET policies are significantly better than random and sometimes even more effective than training a state-of-the-art deep reactive policy from scratch.

## 1 Introduction

A Relational Markov Decision Process (RMDP) (Boutilier, Reiter, and Price 2001) is a first-order, predicate calculus-based representation for expressing instances of a probabilistic planning domain with a possibly unbounded number of objects. An RMDP *domain* has object types, relational state predicate and action symbols that are applied over objects, first order transition templates that specify probabilistic effects associated with action symbols, and a first-order reward structure. A domain *instance* additionally specifies a set of

objects and a start state, thus defining a ground MDP with a known start state (Kolobov, Mausam, and Weld 2012)). *Relational* planners aim to produce a single *generalized* policy that can yield a ground policy for *all* instances of the domain, with little instance-specific computation. *Domain-independent* planners are representation-specific, but domain-agnostic, making them applicable to all domains expressible in the language. In this paper, we design a domain-independent relational planner.

RMDP planners, in their vision, expect to scale to very large problem sizes by exploiting the first-order structures of a domain – thereby reducing the curse of dimensionality. Traditional RMDP planners attempted to find a generalized *first-order* value function or policy using symbolic dynamic programming (Boutilier, Reiter, and Price 2001), or by approximating them via a function over first-order basis functions (e.g., (Guestrin et al. 2003; Sanner and Boutilier 2009)). Unfortunately, these methods met with rather limited success, for e.g., no relational planner participated in International Probabilistic Planning Competition (IPPC)<sup>1</sup> after 2006, even though all competition domains were relational. We believe that this lack of success may be due to the inherent limitations in the representation power of a basis function-based representation. Through this work, we wish to revive the research thread on RMDPs and explore if neural models could be effective in representing these first-order functions.

We present **Symbolic NetWork** (SYMNET), the first domain-independent neural relational planner that computes generalized policies for RMDPs that are expressed in the symbolic representation language of RDDDL (Sanner 2010). SYMNET outputs its generalized policy via a neural model whose all parameters are specific to a domain, but tied among all instances of that domain. So, on a new test instance, the policy can be applied out of the box using pre-trained parameters, i.e., without any retraining on the test instance. SYMNET is domain-independent because it converts an RDDDL domain file (and instance files) completely automatically into neural architectures, without any human intervention.

SYMNET architecture uses two key ideas. First, it visualizes each state of each domain instance as a graph, where

nodes represent the *object tuples* that are valid arguments to some relational predicate. An edge between two nodes indicates that an action causes predicates over these two nodes to interact in the instance. The values of predicates in a state act as features for corresponding nodes. SYMNET then learns node (and state) embeddings for these graphs using graph neural networks. Second, SYMNET learns a neural network to represent the policy and value function over this graph-structured state. To learn these in an instance-independent way, we recognize that most ground actions are a first-order action symbol applied over some object tuple. SYMNET scores such ground actions as a function over the action symbol and the relevant embeddings of object tuples. After training all model parameters using reinforcement learning over training instances of a domain, SYMNET architecture can be applied on any new (possibly larger) test problem without any further retraining.

We perform experiments on nine RDDDL domains from IPPC 2014 (Grzes, Hoey, and Sanner 2014). Since no planner exists that can run without computation on a given instance, we compare SYMNET to random policies (lower bound) and policies trained from scratch on the test instance. We find that SYMNET obtains hugely better rewards than random, and is quite close to the policies trained from scratch – it even outperforms them in 28% instances. Overall, we believe that our work is a step forward for the difficult problem of domain-independent RMDP planning. We release the code of SYMNET for future research.<sup>2</sup>

## 2 Background and Related Work

### 2.1 Probabilistic Planning

**Markov Decision Process (MDP):** A (ground) finite-horizon MDP (Bellman 1957; Puterman 1994) with a known start state is formalized as a tuple  $\langle S, A, T, R, H, s_0, \gamma \rangle$ , where  $S$  is the set of states,  $A$  is the set of actions,  $T$  is the transition model  $S \times A \times S \rightarrow [0, 1]$ ,  $R$  is the reward model  $S \times A \times S \rightarrow \mathbb{R}$ ,  $H$  is the horizon and  $s_0$  is the start state, and  $\gamma$ , the discount factor.

**Relational Markov Decision Process (RMDP):** An RMDP  $\langle \mathcal{C}, \mathcal{SP}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{R}, H, s_0, \gamma \rangle$  is a first-order representation of a factored MDP (Boutilier, Reiter, and Price 2001), expressed via objects, predicates and functions. Here,  $\mathcal{C}$  is a set of classes (types),  $\mathcal{SP}$  is the set of state predicate symbols,  $\mathcal{A}$  is a set of action predicate symbols,  $\mathcal{O}$  represents a set of domain objects, each associated with single type from  $\mathcal{C}$ . It is a first order representation because different sets of objects  $\mathcal{O}$  can construct different ground MDPs.

Each predicate symbol is declared to take as argument a tuple of object types. A predicate symbol (action symbol) applied over a type-consistent tuple of object variables forms a state variable (respectively, ground action). A ground-state  $s$  is, thus, a complete assignment of all predicate symbols  $\mathcal{SP}$  applied on all type-consistent object tuples from  $\mathcal{O}$  (also denoted by  $\mathcal{SP}_{\mathcal{O}}$ ). Similarly, the set of all ground actions ( $\mathcal{A}$ ) can be defined as  $\mathcal{A}_{\mathcal{O}}$  – all-action symbols applied on all type-consistent object tuples. We also denote the ground

state space  $S$  by  $\mathcal{P}(\mathcal{SP}_{\mathcal{O}})$ , where  $\mathcal{P}$  denotes the powerset. Transition and reward models for an RMDP are defined at the schema level through different languages, e.g., PPDDL (Younes et al. 2005) and, our focus, RDDDL (Sanner 2010).

Research in Relational MDPs explores ways to represent and construct first-order (generalized) value functions or policies, which can be used directly on a new test instance. Example representations for these include regression trees (Mausam and Weld 2003), decision lists (Fern, Yoon, and Givan 2006), extensions of algebraic decision diagrams (Joshi and Khardon 2011), and linear combinations of basis functions (Guestrin et al. 2003; Sanner and Boutilier 2005). We believe a reason for limited success of RMDP algorithms is the inherent limitation of these representations. We study the use of deep neural models for representing such generalized functions. To the best of our knowledge, we are the first to develop relational planners for domains expressed in RDDDL.

**Relational Dynamic Decision Language (RDDDL):** RDDDL has been the language of choice for the last three IPPCs. It divides  $\mathcal{SP}$  into non-fluent ( $\mathcal{NF}$ ) and fluent ( $\mathcal{F}$ ) symbols. Non-fluents are the state variables that do not change with time in a given instance, but may be different across problem instances. Fluents represent state variables that change with time (due to actions or natural dynamics). RDDDL splits an RMDP into two separate files, one for the whole domain (that has types, predicates, transitions, and rewards), and the other for the instance (that has objects, non-fluent values, and fluent values for initial state). RDDDL uses additive rewards, total reward is sum of local rewards collected for satisfying different properties in a state. It factors transition function via an underlying DBN semantics. There exist algorithms that convert an RDDDL instance into ground DBN (Sanner 2010).

**Running Example:** We use a simplified Wildfire domain as our example. It has a grid where each cell may have fuel, causing it to burn. The goal is to have the least damage to the grid by either putting out the fire or cutting out the fuel supply. The DBN for the domain is show in Figure 1.

There are two classes,  $\mathcal{C} = \{x_{pos}, y_{pos}\}$ :  $x$  and  $y$  coordinate of the grid cell. Domain has two fluent symbols  $\mathcal{F} = \{burning, out-of-fuel\}$ , representing the current burning state and the fuel state of the cell. Both fluent symbols take a cell  $(x, y)$  as its arguments. The non-fluents represent costs and topology,  $\mathcal{NF} = \{CostTgtBurn, CostNTgtBurn, Neighbour, Target\}$ . The non fluent symbol *Neighbour* takes four arguments  $(x, y, x', y')$ , since it defines the topology of the grid. *Target* has arguments  $(x, y)$ .  $\mathcal{A} = \{put-out, cut-out, finisher\}$ . First two action symbols take arguments  $(x, y)$  – they put out fire and cut out fuel supply at a cell. There is one global action *finisher*, which puts out fire in all the cells simultaneously. Reward (negative) in each time step adds *CostTgtBurn* for each target cell that is burning and *CostNTgtBurn* for each non-target cell that is burning. In a problem instance, say there are three objects  $\mathcal{O} = \{x1, x2, y1\}$ . This implies a problem with two cells  $(x1, y1)$ , and  $(x2, y1)$ . Say the target cell is  $(x1, y1)$  and that these are connected, i.e., *Neighbour* $(x1, y1, x2, y1) = 1$ .

<sup>2</sup><https://github.com/dair-iitd/symnet>

## 2.2 Reinforcement Learning

Reinforcement Learning (RL) refers to planning problems without known transition and rewards, necessitating learning from experience. State of the art approaches for RL are neural, which approximate policy and value functions through deep neural models. We use the Asynchronous Advantage Actor-Critic (A3C) (Mnih et al. 2016) as our underlying RL algorithm. A3C uses two neural networks 1)  $\theta_\pi$  to represent the policy (mapping from a state to distribution over actions) and 2)  $\theta_V$  to estimate the state value (long term discounted reward starting in a state). Policy parameters are optimized to prefer an action that increases a state’s advantage function – difference between its value when taking that action and its overall value. Value parameters optimize the MSE loss between observed and predicted long-term rewards.

## 2.3 Graph Neural Networks

Graph Neural Networks input a graph and learn latent space embeddings for each node, based on the its individual features and the local connectivity structure. Examples include Graph Convolution Networks (GCN) (Kipf and Welling 2017), and Graph Attention Networks (GAT) (Velickovic et al. 2017). We use GAT, which computes a node embedding by using a weighted attention for each of neighbouring nodes. Specifically output node embedding  $\bar{v}_i' = \sigma(\frac{1}{K} \sum_{k=1}^K \sum_{j \in N_i} \alpha_{ij}^k \mathbf{W}^k \bar{v}_j)$ , where  $\bar{v}_i$  is the input feature of node  $v_i$ ,  $N_i$  are its neighbours,  $\mathbf{W}^k$  is a trainable weight matrix,  $k$  is the multi-head hyperparameter and  $\alpha_{ij}^k = \text{softmax}_j(f(\mathbf{W}^k \bar{v}_i, \mathbf{W}^k \bar{v}_j))$  is the normalized self attention coefficient for any non-linear function ( $f$ ), which in our implementation is LeakyReLU (Xu et al. 2015).

## 2.4 Transfer Learning for Probabilistic Planning

There having been several classical (Taylor and Stone 2009; Sorg and Singh 2009; Atkeson and Schaal 1997) and neural (Parisotto, Ba, and Salakhutdinov 2015; Matiisen et al. 2017; Garnelo, Arulkumaran, and Shanahan 2016; Higgins et al. 2017) approaches for transfer learning in RL. Recent work studies transfer learning for symbolic planning problems, e.g., Groshev et al. (2018) for deterministic planning problems. ASNets, Action-Schema Networks (Toyer et al. 2018; Shen et al. 2019), tackle a problem similar to ours but for goal-oriented subset of PPDDL. While an RDDDL instance can be converted automatically into propositional PPDDL, an RDDDL domain cannot always be converted into relational PPDDL – hence we cannot directly compare against ASNets. Issakkimuthu et al. (2018) devise a neural framework to learn a policy for (ground) RDDDL MDPs from scratch. Their constraint on non-transferability is due to the fixed size of fully connected layers in the neural network. TORPIDO achieves transfer across RDDDL problem instances of the same domain (Bajpai, Garg, and Mausam 2018); it can only transfer over *equi-sized* problems due to its fixed size decoder.

Our previous work TRAPSNET is closest to SYMNET, as it can transfer to different-sized instances of an RMDP (Garg, Bajpai, and Mausam 2019). It constructs a graph, which uses single object as nodes, and non-fluent based edge. It encodes each node in embedding space and computes the score for a

ground action based on the applied action template, and object embedding. However, TRAPSNET makes the restrictive assumptions that the domains have exactly one binary non-fluent, and all the rest are unary fluents or non-fluents, and each action symbol is parameterized by *exactly* one object. These assumptions don’t hold in several RDDDL domains.

## 3 Problem Formulation

Given an RMDP domain  $D = \langle \mathcal{C}, \mathcal{SP}, \mathcal{A}, \mathcal{T}, \mathcal{R}, H, \gamma \rangle$  expressed in RDDDL, we wish to learn a generalized policy  $\pi^D$ , which can be applied to all instances of  $D$  and maximizes the discounted sum of expected rewards over a finite horizon  $H$ . Given a test problem instance  $I_t = \langle O, s_0 \rangle$  from  $D$ , this generalized policy can yield an instance-specific policy  $\pi^D(I_t) : \mathcal{P}(\mathcal{SP}_O) \rightarrow \mathcal{A}_O$ , without any training on  $I_t$ . The RMDP learning problem can be seen in terms of multi-task learning over several problem instances in  $D$ : given  $N$  randomly selected problem instances  $I_1, I_2, \dots, I_N$  (possibly of different sizes) from  $D$ , we wish to learn the weights  $\phi$  of a neural network, such that  $\pi^D(I_i; \phi)$  is a good (high-reward) policy for problem instance  $I_i$ . A good generalized policy is one which, without training, achieves high reward values on the new instance  $I_t$ .

## 4 The SYMNET Framework

We now present SYMNET’s architecture for training a generalized policy for a given RMDP domain. We follow existing research to hypothesize that for any instance of a domain, we can learn a representation of the current state in a latent space and then output a policy in latent space, which is decoded into a ground action. To achieve this, SYMNET uses three modules: (1) problem representation, which constructs an instance graph for every problem instance, (2) representation learning, which learns embeddings for every node in the instance graph, and for the state, and (3) policy decoder, which computes a value for every ground action, outputting a mixed policy for a given state. All parameters of representation learning and policy learning modules are shared across all instances of a domain. SYMNET’s full architecture is shown in Figure 2.

### 4.1 Problem Representation

We follow TRAPSNET, in that we continue the general idea of converting an instance into an instance graph and then learning a graph encoder to handle different-sized domains. However, the main challenge for a general RMDP, one that does not satisfy the restricted assumptions of TRAPSNET, is in defining a coherent graph structure for an instance. The first key question is what should be a node in the instance graph. TRAPSNET’s approach was to use a single object as nodes, as all fluents (and actions) in its domains took single objects as arguments. This may not work for a general RMDP since it’s fluents and actions may take several objects as arguments. Secondly, how should edges be defined. Edges represent the interaction between nodes. TRAPSNET defined them based on the one binary non-fluent in its domain. A general RMDP may not have any non-fluent symbol or may have many (possibly higher-order) non-fluents.

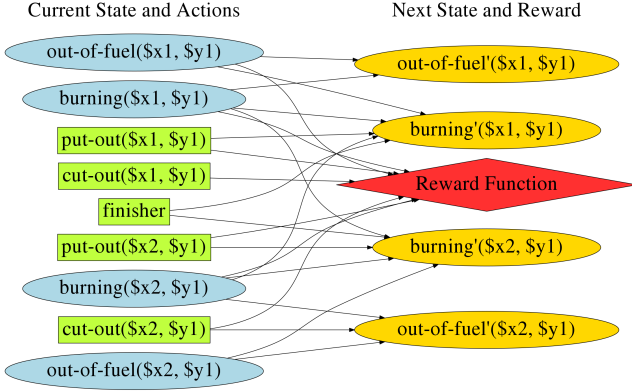


Figure 1: DBN for a modified wildfire problem.

Last but not least, the real domain-independence for SYMNET can be achieved only when it parses an RDDDL domain file without any human intervention. This leads to a novel challenge of reconciling multiple different ways in RDDDL to express the same domain. In our running example, connectivity structure between cells may be defined using non-fluents  $y$ -neighbour( $y, y'$ ),  $x$ -neighbour( $x, x'$ ), or using a quaternary non-fluent  $neighbour(x, y, x, y')$ . Since both these representations represent the same problem, an ideal desideratum is that the graph construction algorithm leads to the same instance graph in both cases. But, this is a challenge since the corresponding RDDDL domains may look very different. While, in general, this problem seems too hard to solve, since it is trying to judge logical equivalence of two domains, SYMNET attempts to achieve the same instance graphs in case the equivalence is within non-fluents.

To solve these problems, we make the observation that dynamics of an RDDDL instance ultimately compile to a ground DBN with nodes as state variables (fluent symbols applied on object tuples) and actions (action symbols applied on object tuples).<sup>3</sup> DBN exposes a connectivity structure that determines which state variables and actions directly affect another state variable. It additionally has conditional probability tables (CPTs) for each transition. Figure 1 shows an example of a DBN for our running example instance. Here, left column is for current time step, and right for the next one. The edges represent which state and action variables affect the next state-variable. We note that the ground DBN does not expose non-fluents since its values are fixed, and their dependence can be compiled directly into CPTs.

SYMNET converts a ground DBN to an instance graph. It constructs a node for every unique *object tuple* that appears as an argument in any state variable in the DBN. Moreover, two nodes are connected if the state variables associated with two nodes influence each other in the DBN through some action. This satisfies all our challenges. First, it goes beyond an object as a node, but only defines those nodes that are likely important in the instance. Second, it defines a clear semantics of edges, while maintaining its intuition

<sup>3</sup>done automatically using code from <https://github.com/ssanner/rddlsim>

of “directly influences.” Finally, it can handles some variety of non-fluent representations for the same domain. Since the DBN does not even expose non-fluent state variables, and compiles them away, the same instances encoded with different non-fluent representations often yield same ground DBNs and thus the same instance graphs.

**Construction of Instance Graph:** We now formally describe the conversion of a DBN into a directed instance graph,  $G = (V, E)$ , where  $V$  is the set of vertices and  $E$  is the set of edges.  $G$  is composed of  $\mathcal{K} = |\mathcal{A}| + 1$  disjoint subgraphs  $G_j = (V_j, E_j)$ . Intuitively, each graph  $G_j$  has information about influence of each individual action symbol  $a_j \in \mathcal{A}$ .  $G_{\mathcal{K}}$  represents the influence of the full set  $\mathcal{A}$ , and also the natural dynamics. In our example  $\mathcal{K} = 4$  since we have three action symbols: put-out, cut-out and finisher.

To describe the formal process, we define three analogous sets:  $O_f$ ,  $O_{nf}$  and  $O_a$ .  $O_f$  represents the set of all object tuples that act as a valid argument for any fluent symbol.  $O_{nf}$  and  $O_a$  are analogous sets for non-fluent and action symbols. In our running example,  $O_f = \{(x1, y1), (x2, y1)\}$ ,  $O_{nf} = \{(x1, y1), (x2, y1), (x1, y1, x2, y1), (x2, y1, x1, y1)\}$ , and  $O_a = \{(x1, y1), (x2, y1)\}$ . Nodes in the instance graph associate with object tuples. We use  $o_v$  to denote the object tuple associated with node  $v$ . SYMNET converts a DBN into an instance graph as follows:

1. The distinct object tuples in fluents form the nodes of the graph, i.e.  $V_j = \{v | o_v \in O_f\}, \forall j$ . For the example, each  $V_j =$  different copies of  $\{(x1, y1), (x2, y1)\}$ .
2. We add an edge between two nodes in  $G_j$  if some state variables corresponding to them are connected in the DBN through  $a_j$ . Formally,  $E_j(u, v) = 1$ , if  $\exists f, g \in F, \exists o_a \in O_a, j \in \{1, \dots, |\mathcal{A}|\}$  s.t. the transition dynamics ( $T^f$ ) for state variable  $g'(o_v)$  and action  $a_j(o_a)$  depend on state variable  $f(o_u)$  or  $f'(o_u)$ . For the running example, there is no edge between  $(x1, y1)$  and  $(x2, y1)$  since cut-out, put-out or finisher’s effects on one cell do not depend on any other cell.
3. We add an edge between two nodes in  $G_{\mathcal{K}}$  if some state variables corresponding to them are connected in the DBN (possibly through natural dynamics). I.e.,  $E_{\mathcal{K}}(u, v) = 1$ , if  $\exists f, g \in F$  s.t. there is an edge from  $f(o_u)$  (or  $f'(o_u)$ ) and  $g'(o_v)$  in the DBN. For the example,  $E_4((x1, y1), (x2, y1)) = 1$  as there is an edge between  $burning(x1, y1)$  and  $burning'(x2, y1)$  since fire propagates to neighboring cells through natural dynamics. Similarly,  $E_4((x2, y1), (x1, y1)) = 1$ .
4. As every node influences itself, self loops are added on each node.  $E(v, v) = 1, \forall v \in V$ .

For each node  $v \in V$ , we additionally construct a feature vector ( $h(v)$ ) which consists of fluent feature vector ( $h^f(v)$ ) and non-fluent feature vector ( $h^{nf}(v)$ ), such that  $h = concat(h^f, h^{nf})$ . The feature vector for all nodes for the same object tuple is the same. The feature vector is constructed as follows:

1. The fluent features for each node are obtained from the state of the problem instance. The values of state variables corresponding to a node are added as feature to that node. Whenever a fluent symbol cannot take a node as an argument, we add zero as the feature for it. Formally,

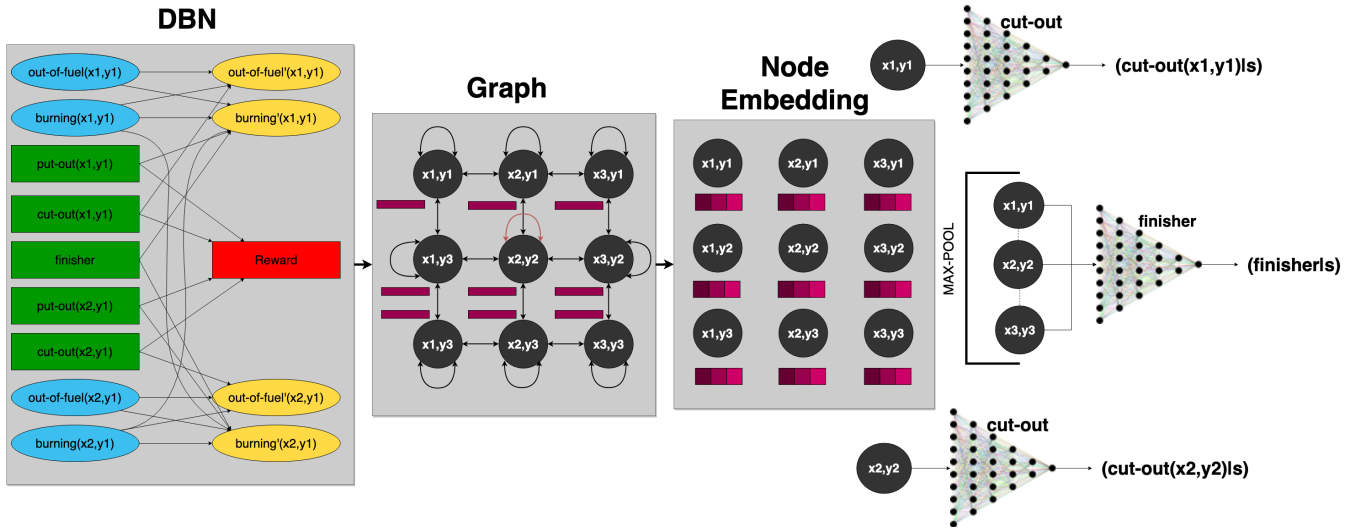


Figure 2: Policy network for SYMNET demonstrated on  $2 \times 1$  wildfire domain. Fully Connected Network is used in Action Decoder.

$h^f(v)_i = g_i(o_v)$  if  $g_i \in \mathcal{F}, v \in V$  and  $o_v$  is an argument of  $g_i$ , otherwise,  $h^f(v)_i = 0, \forall i = 1 \dots |\mathcal{F}|$ . For the running example, we have two state-fluents. Hence,  $h^f((x1, y1)) = [burning(x1, y1), out-of-fuel(x1, y1)]$ .

- The non-fluent feature vector for each node is obtained from the RDDDL file. The values of non-fluents defined on the node, and additionally any unary non-fluents where the argument intersects the node are added as the features for the node. The default value is obtained from the domain file while the specific value (if available) is obtained from the instance file. Formally,  $h^{nf}(v)_i = g_i(o_{nf})$  if  $g_i \in \mathcal{NF}, v \in V, o_{nf} \in O_{nf}, ((o_v = o_{nf}) \vee (|o_{nf}| = 1 \wedge o_{nf} \subset o_v))$ , otherwise,  $h^{nf}(v)_i = 0, \forall i = 1 \dots |\mathcal{NF}|$ . In our example,  $h^{nf}((x1, y1)) = [target(x1, y1)]$ .

We note that the size of feature vector on each node depends on the domain, but is independent of the number of objects in the instance – there are a constant number of feature values per state predicate symbol. This allows variable-sized instances of the same domain to use the same representation.

## 4.2 Representation Learning

SYMNET runs GAT on the instance graph to obtain node embeddings  $\bar{v}$  for each node  $v \in V$ . It then constructs tuple embedding for each object tuple by concatenating node embeddings of all associated nodes. Formally, let  $O_V = \{o_v | v \in V\}$ . For  $o \in O_V$ , the tuple embedding  $\bar{o} = \text{concat}(\bar{v})$ , over all  $v$  s.t.  $o_v = o$ . SYMNET also computes a state embedding  $\bar{s}$  by taking a dimension-wise max over all tuple embeddings, i.e.,  $\bar{s} = \text{MaxPool}_{o \in O_V}(\bar{o})$ .

## 4.3 Policy Decoder

SYMNET maps latent representations  $\bar{o}$  and  $\bar{s}$  into a state value  $V(s)$  (long-term expected discounted reward starting in state  $s$ ) and mixed policy  $\pi(s)$  (probability distribution

over all ground actions). This is done using a value decoder and a policy decoder, respectively.

There are several challenges in designing a (generalized) policy decoder. First, the action symbols may take multiple objects as arguments. Second, and more importantly, action symbols may even take those object tuples as arguments that do not correspond to any node in the instance graph. This will happen if an object tuple (in  $O_a$ ) is not an argument to any fluent symbol, i.e.,  $\exists o_a$  s.t.  $o_a \in O_a \wedge o_a \notin O_f$ . Adding these object tuples as nodes in the instance graph may not work, since we won't have any natural features for those nodes.

In response, we design a novel framework for policy and value decoders. The decoders consist of fully connected layers, the input to which are a subset of the tuple embeddings  $\bar{o}$ . SYMNET uses the following rules to construct decoders:

- The number of decoders is constant for a given domain and is equal to the number of distinct action symbols ( $|\mathcal{A}|$ ). For the running example, three different decoders for each policy and value decoding are constructed, namely *cut-out*, *put-out* and *finisher*.
- The input to a decoder is the state embedding  $\bar{s}$  concatenated with embeddings of object tuples corresponding to the state variables affected by the action in the DBN. In running example, *put-out*( $x1, y1$ ) action takes only the tuple embedding of ( $x1, y1$ ) as input. However, the number of state-variables being affected by a ground action might vary across instances of the same domain. For example, the *finisher* action affects all cells. To alleviate this, we use size-independent max pool aggregation over the embeddings of all affected tuple embeddings to create a fixed-sized input.
- Decoder parameters are specific to action symbols and not to ground actions. In running example, *put-out*( $x1, y1$ ) will be scored using embedding of ( $x1, y1$ ); similarly, for ( $x2, y1$ ). But, both scorings will use a single parameter set specific to *put-out*.

Table 1:  $\alpha_{symnet}(0)$  values of SYMNET. Bold values represent over 90% the score of max performance.

Instance	5	6	7	8	9	10	
Domain	AA	<b>0.93 ± 0.01</b>	<b>0.94 ± 0.01</b>	<b>0.94 ± 0.01</b>	<b>0.92 ± 0.02</b>	<b>0.95 ± 0.03</b>	<b>0.91 ± 0.05</b>
	CT	0.87 ± 0.16	0.78 ± 0.14	<b>1.00 ± 0.07</b>	<b>0.98 ± 0.13</b>	<b>0.99 ± 0.04</b>	<b>1.00 ± 0.05</b>
	GOL	<b>0.96 ± 0.06</b>	<b>1.00 ± 0.05</b>	0.65 ± 0.05	0.83 ± 0.03	<b>0.95 ± 0.04</b>	0.64 ± 0.08
	Nav	<b>0.99 ± 0.01</b>	<b>1.00 ± 0.01</b>	<b>0.99 ± 0.01</b>	<b>1.00 ± 0.01</b>	<b>1.00 ± 0.02</b>	<b>1.00 ± 0.02</b>
	ST	<b>0.91 ± 0.05</b>	0.84 ± 0.02	0.86 ± 0.05	0.85 ± 0.05	0.81 ± 0.02	0.89 ± 0.03
	Sys	<b>0.96 ± 0.03</b>	<b>0.98 ± 0.02</b>	<b>0.98 ± 0.02</b>	<b>0.97 ± 0.02</b>	<b>0.99 ± 0.01</b>	<b>0.96 ± 0.03</b>
	Tam	<b>0.92 ± 0.07</b>	<b>1.00 ± 0.12</b>	<b>0.98 ± 0.06</b>	<b>1.00 ± 0.12</b>	<b>1.00 ± 0.12</b>	<b>0.95 ± 0.06</b>
	Tra	0.85 ± 0.18	<b>0.93 ± 0.06</b>	0.88 ± 0.21	0.74 ± 0.17	<b>0.94 ± 0.12</b>	0.87 ± 0.13
	Wild	<b>0.99 ± 0.01</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.01</b>	<b>1.00 ± 0.01</b>

- The policy decoder computes scores of all ground actions, which are normalized using softmax to output the final policy in a state. For  $I_t$ , the highest probability action is selected as the final action.
- All value outputs are summed to give the final value for that state. This modeling choice reflects the additive reward aspect of many RDDDL domains.

#### 4.4 Learning

While construction of SYMNET architecture is heavily dependent on the RDDDL domain and instance files, actual training is done via model-free reinforcement learning approach of A3C (Mnih et al. 2016). RL learns from interactions with environment – SYMNET simulates the environment using RDDDL-specified dynamics. Use of model-based planning algorithms for this purpose is left as future work. We formulate training of SYMNET as a multi-task learning problem (see Section 3), so that it generalizes well and does not overfit on any one problem instance. The parameters for the state encoder, policy decoder, and value decoder are learned using updates similar to that in A3C. SYMNET’s loss function for the policy and value network is the same as that in the A3C paper (summed over the multi-task problem instances).

As constructed, SYMNET’s number of parameters is independent of the size of the problem instance. Hence, the same network can be used for problem instances of any size. After the learning is completed, the network represents a generalized policy (or value), since it can be directly used on a new problem instance to compute policy in a single forward pass.

## 5 Experiments

Our goal is to estimate the effectiveness of SYMNET out-of-the-box policy for a new problem in a domain. Unfortunately, there are no available transfer algorithms for general RDDDL RMDPs. So, we first compare it against a random policy, because that is the best we can do currently with no time to train. To further understand the overall quality of the generalized policy, we also compare it against several neural models that train from scratch on the test instance. We also compare it against state-of-the-art online planner PROST (Keller and Eyerich 2012).

### 5.1 Domains and Experimental Setting

We show all our results on nine RDDDL domains used IPPC 2014: Academic Advising (AA), Crossing Traffic (CT), Game of Life (GOL), Navigation (NAV), Skill Teaching (ST), Sysadmin (Sys), Tamarisk (Tam), Traffic (Tra), and Wildfire (Wild). We describe the domains, number of state fluents & non-fluents, action fluents in supplementary material. The RL agent is trained to learn the generalized policy on smaller sized instances. We use IPPC problem instances 1, 2, and 3 of each domain for the multi-task training of SYMNET network. A3C uses a discounted parameter of 0.99. In the spirit of domain-independent planning, we use the *same* hyperparameters for each domain. The embedding module for GAT uses a neighborhood of 1 and an output feature size of 6. We then use a fully connected layer of output 20 dimensions to get an embedding from each of the tuple embedding outputs by GAT. All layers use a leaky ReLU activation and a learning rate of  $10^{-3}$ . We train the network using RMSProp (Ruder 2016) on a single Nvidia K40 GPU. SYMNET is trained for each domain for twelve hours (4 hours for each instance).

### 5.2 Comparison Algorithms and Metrics

As there does not exist any previous method for learning over Relational RDDDL MDPs, we can only compare against a random policy. However, this experiment can only show the difference from a random policy, but cannot evaluate the overall goodness of the generalized policy. For that, we compare against several (potentially upper bound) policies that are not directly comparable to SYMNET in their experimental settings. For our first such experiment, we use TORPIDO as the state-of-the-art deep reactive policy. Note that we do not use their transfer method, but train the network from scratch on the problem instance. This is because it can only transfer across equi-sized instances. Still, it is an upper bound as TORPIDO trained on the test instance is compared against SYMNET trained on other smaller instances, but not the test instance. Similarly, we also compare against SYMNET architecture itself, trained from scratch on the test instance (named SYMNET-s). The main difference between TORPIDO and SYMNET architectures is that TORPIDO has a much higher capacity since it models each ground action explicitly. On three domains where TRAPSNET is applicable, we also compare against TRAPSNET policies out of the box. Finally, we also compare against the state-of-the-art *online* planner,

Table 2: Comparison of SYMNET against SYMNET-s (SYM) architecture trained from scratch and TORPIDO (TOR) architecture trained from scratch. We compare out-of-the-box SYMNET to others after 12 hours of training. INF is used when SYM or TOR achieved minimum possible reward and hence SYMNET was infinitely better.

Domain	SYM	TOR	Domain	SYM	TOR	Domain	SYM	TOR	Domain	SYM	TOR
AA 5	1.09	0.99	GOL 5	<b>1.35</b>	<b>1.49</b>	ST 5	1.11	0.94	Tam 5	<b>INF</b>	<b>2.33</b>
AA 6	1.78	0.95	GOL 6	<b>1.57</b>	<b>1.69</b>	ST 6	1.21	0.90	Tam 6	<b>27.71</b>	<b>8.13</b>
AA 7	1.21	0.98	GOL 7	1.08	0.76	ST 7	1.10	0.87	Tam 7	<b>17.81</b>	<b>4.83</b>
AA 8	1.31	0.97	GOL 8	2.22	0.87	ST 8	1.14	0.90	Tam 8	<b>2.74</b>	<b>15.56</b>
AA 9	1.39	0.95	GOL 9	<b>1.86</b>	<b>1.31</b>	ST 9	1.13	0.81	Tam 9	<b>24.94</b>	<b>13.07</b>
AA 10	1.32	0.93	GOL 10	1.25	0.68	ST 10	1.30	0.95	Tam 10	<b>2.35</b>	<b>7.99</b>
CT 5	<b>1.34</b>	<b>1.39</b>	Nav 5	<b>10.84</b>	<b>INF</b>	Sys 5	<b>1.03</b>	<b>2.89</b>	Tra 5	1.78	0.86
CT 6	<b>INF</b>	<b>1.56</b>	Nav 6	<b>INF</b>	<b>INF</b>	Sys 6	<b>1.33</b>	<b>1.20</b>	Tra 6	<b>1.56</b>	<b>1.39</b>
CT 7	<b>1.13</b>	<b>1.12</b>	Nav 7	<b>INF</b>	<b>INF</b>	Sys 7	<b>1.56</b>	<b>2.45</b>	Tra 7	<b>3.28</b>	<b>1.13</b>
CT 8	<b>1.55</b>	<b>1.23</b>	Nav 8	<b>INF</b>	<b>INF</b>	Sys 8	<b>1.46</b>	<b>1.60</b>	Tra 8	1.13	0.81
CT 9	<b>1.35</b>	<b>1.16</b>	Nav 9	<b>INF</b>	<b>INF</b>	Sys 9	<b>1.38</b>	<b>1.17</b>	Tra 9	<b>2.50</b>	<b>1.08</b>
CT 10	<b>1.22</b>	<b>4.99</b>	Nav 10	<b>INF</b>	<b>INF</b>	Sys 10	<b>1.18</b>	<b>1.50</b>	Tra 10	<b>1.53</b>	<b>1.86</b>
Wild 5	<b>1.03</b>	<b>1.13</b>	Wild 7	<b>1.03</b>	<b>1.13</b>	Wild 9	<b>1.01</b>	<b>13.14</b>			
Wild 6	<b>1.01</b>	<b>1.01</b>	Wild 8	<b>1.00</b>	<b>1.09</b>	Wild 10	<b>34.80</b>	<b>11.19</b>			

Table 3: Comparison of TRAPSNET with SYMNET on three domains as published in (Garg, Bajpai, and Mausam 2019). Label: AA - Academic Advising, GOL - Game Of Life, Sys - Sysadmin

Instance	5	6	7	8	9	10
Domain AA	<b>1.12</b>	<b>1.17</b>	<b>1.12</b>	<b>1.27</b>	<b>1.26</b>	<b>1.40</b>
Domain GOL	0.96	<b>1.04</b>	0.69	<b>1.00</b>	0.97	<b>1.50</b>
Domain Sys	<b>1.01</b>	<b>1.55</b>	<b>1.33</b>	<b>1.39</b>	<b>1.21</b>	<b>1.17</b>

### PROST.

After training algorithm  $alg$  for  $t$  hours, we simulate its output policy 200 times (for  $H$  steps each) from the start state. We average the discounted rewards to estimate the expected long term discounted reward of that policy, denoted by  $V_{alg}(t)$ . To be able to compare across domains and problems and reward ranges, we report a normalized metric  $\alpha_{alg}(t) = \frac{V_{alg}(t) - V_{min}}{V_{max} - V_{min}}$  where  $V_{min}$  and  $V_{max}$  are the minimum, and the maximum expected discounted rewards obtained at any time by any of the four comparison algorithms on a given instance. This number lies between 0 and 1, with 1 being the best-found reward, and 0 being the random policy’s reward. All algorithms are trained independently 5 times and the average result is reported. During training from scratch, all networks start with a random policy and hence have their  $\alpha(0)$  values as 0. However, that is not true for SYMNET as it is pre-trained on the domain. To compare against other training approaches directly, we compute  $\beta_{alg}(t) = \frac{\alpha_{symnet}(0)}{\alpha_{alg}(t)}$ . A value higher than 1 suggests that SYMNET out-of-the-box outperformed  $alg$  trained for  $t$  hours, and less than 1 implies SYMNET performed worse.

## 5.3 Results

**Comparison against Random Policy:** We report the val-

ues of  $\alpha_{symnet}(0)$  in Table 1. Since the random policy is 0, we notice that on all six problem instances from the nine domains, SYMNET performs enormously better than random. We highlight the instances where our method achieves over 90% of the max reward obtained by any algorithm for that instance. We see that SYMNET with no training achieves over 90% the max reward on 40 instances and over 80% in 50 out of 54 instances. We also show that our method performs the best out-of-the-box in 28 instances. This is our main result, and it highlights that SYMNET takes a major leap towards the goal of computing generalized policies for the whole RMDP domain, and can work on a new instance out of the box.

**Comparison against Training from Scratch:** We now compare SYMNET against the expected discounted rewards obtained by TORPIDO and SYMNET-s, when they are trained from scratch for 12 hours on the test problem. We note that these numbers are not directly comparable, since in one case, the model has been trained on other instances of the domain, but not trained on the test problem at all, and in the other case the models are trained from scratch on the test. That said, this comparison is likely a good indicator of the absolute performance of SYMNET.

Table 2 reports the values for  $\beta_{torpido}(12)$  and  $\beta_{symnet-s}(12)$ . We notice that, surprisingly, SYMNET policy with no training is better than both methods on several instances. Against SYMNET trained from scratch, it is better on all instances, although its edge over TORPIDO is limited to 37 out of 54. We hypothesize that this excellent performance is due to the multi-task learning aspect of SYMNET, where it is able to reach some generalized policy of a domain that is not found on the specific instance even after training for 12 hours.

In 17 out of 54 instances, SYMNET lags behind TORPIDO, which is not surprising, since TORPIDO has much higher capacity, as discussed earlier. We also notice that the performance of TORPIDO is no better than random for

Table 4: Comparison of PROST with SYMNET. INF is used when PROST returned a policy equal to or worse than a random policy.

	Domain	AA	CT	GOL	Nav	ST	Sys	Tam	Tra	Wild
Instance	5	<b>2.13</b>	0.76	0.48	<b>1.16</b>	0.95	<b>1.13</b>	0.61	0.60	<b>1.39</b>
	6	<b>2.14</b>	0.44	0.57	<b>1.87</b>	0.86	<b>1.24</b>	0.96	0.65	<b>INF</b>
	7	<b>2.18</b>	0.62	0.33	<b>6.42</b>	0.86	<b>1.13</b>	0.70	0.61	<b>INF</b>
	8	<b>1.79</b>	0.37	0.39	<b>45.46</b>	0.90	<b>1.50</b>	0.79	0.51	<b>INF</b>
	9	<b>1.46</b>	0.74	0.44	<b>101.23</b>	0.78	<b>1.21</b>	0.83	0.75	<b>INF</b>
	10	<b>1.46</b>	0.37	0.30	<b>INF</b>	0.93	<b>1.42</b>	0.84	0.64	<b>1.49</b>

Navigation. We attribute this to the sparse and late reward obtained in large instances of this domain, which makes it difficult for TORPIDO to learn a good policy. Because of the late rewards, TORPIDO is not able to reach the goal state at all in 12 hours of training, and hence is not able to improve on the random policy. SYMNET trains well on small instances where path to goal is short and generalizes well. In GOL, SYMNET performs worse, because the nature of policy changes significantly in large instances (e.g. requiring new patterns to survive) which cannot be learned in smaller instances at all.

**Comparison against TRAPSNET:** While TRAPSNET is not applicable in many RMDPs, still, we can compare it with SYMNET on some domains. We compare these on three domains that follow the unary fluents and binary non-fluents constraint: Academic Advising, Game of Life, and SysAdmin. We report  $\beta_{trapsnet}(0)$  in Table 3. It shows that SYMNET outperforms TRAPSNET on 15 out of 18 instances, is comparable on 2 instances and worse on 1 instance. We attribute the success of SYMNET over TRAPSNET to the action-symbol specific graphs ( $G_j$ ), which likely help learn better action dependencies in the embeddings.

**Comparison against ASNets:** Even after significant efforts, we were not able to compare against ASNets, which solves a similar problem for PPDDL domains. Converting an RDDDL domain to PPDDL enumerates all the ground state-variables and loses the RMDP structure. This leads to different domain files for different instances for the same problem domain, due to which ASNets is unable to train. We also tried writing a domain file manually for a few domains, but were not successful due to the unavailability of floating non-fluent values, and due to non-additive reward structure in PPDDL.

**Comparison against PROST:** Finally, we compare against PROST. PROST is a state-of-the-art *online* planner, i.e., it performs interleaved planning and execution, as it builds a new search tree before taking every action, based on the specific state reached. On the other hand, SYMNET outputs an offline policy, which does not need much computation for deciding the next action. Offline and online policies are two very different settings, and these results are *not directly comparable*. Nonetheless, we report  $\frac{V_{symnet}(0) - V_{min}}{V_{prost} - V_{min}}$ . The code of PROST is obtained from the official repository and we use its default settings for this comparison.<sup>4</sup>

We compare our policy with PROST on all 9 domains,

shown in Table 4. We see that on four domains SYMNET achieves a much better performance than PROST. This is rather surprising to us that even after substantial lookahead from the current state, PROST is still not able to compute a good policy. For example, in both Navigation and Wildfire, the rewards are sparse and distant, and PROST is often unable to reach the goal in its planning horizon. In other five domains, PROST is substantially better than SYMNET. This suggests that SYMNET policies are not close to optimal, and further research is needed for making them even stronger. This also points to the future possibility of applying a combination of SYMNET and PROST for the offline setting, not unlike the use of Monte-Carlo Tree Search with deep neural networks in AlphaGo (Silver et al. 2016).

Overall, we find that SYMNET’s generalized policies out-of-the-box are enormously better than random, and can frequently beat other deep neural models trained from scratch on the test instance. However, comparison with PROST suggests that SYMNET policies are not close to optimal and further research is needed to make them even better.

## 6 Conclusion and Future Work

We present the first neural-method for obtaining a generalized policy for Relational MDPs represented in RDDDL. Our method, named SYMNET, converts an RDDDL problem instance into an instance graph, on which a graph neural network computes state embeddings and embeddings for important object tuples. These are then decoded into scores for each ground action. All parameters are tied and size-invariant such that the same model can work on problems of varying sizes. In our experiments, we train SYMNET on small problems of a domain and test them on larger problems to find that they out-of-the-box perform hugely better than random. Even when compared against training deep reactive policies from scratch, SYMNET without training perform better or at par in over half the problem instances.

Our work is an attempt to revive the thread on Relational MDPs and the attractive vision of generalized policies for a domain. However, ours is only one of the first steps. Further investigation is needed to assess how far are SYMNET’s generalized policies from optimal. We strongly believe that there may be even better architectures that could learn near-optimal generalized policies, and the need for retraining or interleaving planning and execution could be rendered unnecessary. We release all our software for use by the research community at <https://github.com/dair-iitd/symnet>.

<sup>4</sup><https://github.com/prost-planner/prost>



## Acknowledgements

We would like to thank Scott Sanner for an extremely insightful discussion on DBNs, which enabled us to work out a solution to convert an RMDP into a graph using DBNs. We also thank Vishal Sharma, Gobind Singh, Parag Singla and the anonymous reviewers for their comments on various drafts of the paper. This work is supported by grants from Google, Bloomberg, IBM and IMG, Jai Gupta Chair Fellowship, and a Visvesvaraya faculty award by Govt. of India. We thank Microsoft Azure sponsorships, and the IIT Delhi HPC facility for computational resources.

## References

- Atkeson, C. G., and Schaal, S. 1997. Robot learning from demonstration. In *Proceedings of the Fourteenth International Conference on Machine Learning (ICML 1997)*, Nashville, Tennessee, USA, July 8-12, 1997, 12–20.
- Bajpai, A.; Garg, S.; and Mausam. 2018. Transfer of deep reactive policies for mdp planning. In Bengio, S.; Wallach, H.; Larochelle, H.; Grauman, K.; Cesa-Bianchi, N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc. 10965–10975.
- Bellman, R. 1957. A Markovian Decision Process. *Indiana University Mathematics Journal*.
- Boutillier, C.; Reiter, R.; and Price, B. 2001. Symbolic dynamic programming for first-order mdps. In *Proceedings of the Seventeenth International Joint Conference on Artificial Intelligence, IJCAI 2001*, Seattle, Washington, USA, August 4-10, 2001, 690–700.
- Fern, A.; Yoon, S. W.; and Givan, R. 2006. Approximate policy iteration with a policy language bias: Solving relational markov decision processes. *J. Artif. Intell. Res.* 25:75–118.
- Garg, S.; Bajpai, A.; and Mausam. 2019. Size independent neural transfer for rddl planning. In *Proceedings of the International Conference on Automated Planning and Scheduling*, 631–636.
- Garnelo, M.; Arulkumaran, K.; and Shanahan, M. 2016. Towards deep symbolic reinforcement learning. *CoRR* abs/1609.05518.
- Groshev, E.; Tamar, A.; Goldstein, M.; Srivastava, S.; and Abbeel, P. 2018. Learning generalized reactive policies using deep neural networks. In *ICAPS*.
- Grzes, M.; Hoey, J.; and Sanner, S. 2014. International Probabilistic Planning Competition (IPPC) 2014. In *ICAPS*.
- Guestrin, C.; Koller, D.; Gearhart, C.; and Kanodia, N. 2003. Generalizing plans to new environments in relational mdps. In *IJCAI*, 1003–1010.
- Higgins, I.; Pal, A.; Rusu, A. A.; Matthey, L.; Burgess, C.; Pritzel, A.; Botvinick, M.; Blundell, C.; and Lerchner, A. 2017. DARLA: improving zero-shot transfer in reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*, Sydney, NSW, Australia, 6-11 August 2017, 1480–1490.
- Issakkimuthu, M.; Fern, A.; and Tadepalli, P. 2018. Training deep reactive policies for probabilistic planning problems. In *ICAPS*.
- Joshi, S., and Khardon, R. 2011. Probabilistic relational planning with first order decision diagrams. *Journal of Artificial Intelligence Research* 41:231–266.
- Keller, T., and Eyerich, P. 2012. PROST: probabilistic planning based on UCT. In *Proceedings of the Twenty-Second International Conference on Automated Planning and Scheduling, ICAPS 2012*, Atibaia, São Paulo, Brazil, June 25-19, 2012.
- Kipf, T. N., and Welling, M. 2017. Semi-supervised classification with graph convolutional networks. In *ICLR*.
- Kolobov, A.; Mausam; and Weld, D. S. 2012. A theory of goal-oriented mdps with dead ends. In *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence, Catalina Island, CA, USA, August 14-18, 2012*, 438–447.
- Mattiisen, T.; Oliver, A.; Cohen, T.; and Schulman, J. 2017. Teacher-student curriculum learning. *CoRR* abs/1707.00183.
- Mausam, and Weld, D. S. 2003. Solving relational MDPs with first-order machine learning. In *ICAPS'03 Workshop on Planning under Uncertainty and Incomplete Information*.
- Mnih, V.; Badia, A. P.; Mirza, M.; Graves, A.; Lillicrap, T. P.; Harley, T.; Silver, D.; and Kavukcuoglu, K. 2016. Asynchronous methods for deep reinforcement learning. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016*, New York City, NY, USA, June 19-24, 2016, 1928–1937.
- Parisotto, E.; Ba, L. J.; and Salakhutdinov, R. 2015. Actor-mimic: Deep multitask and transfer reinforcement learning. *CoRR* abs/1511.06342.
- Puterman, M. 1994. *Markov Decision Processes*. John Wiley & Sons, Inc.
- Ruder, S. 2016. An overview of gradient descent optimization algorithms.
- Sanner, S., and Boutilier, C. 2005. Approximate linear programming for first-order mdps. In *UAI '05, Proceedings of the 21st Conference in Uncertainty in Artificial Intelligence, Edinburgh, Scotland, July 26-29, 2005*, 509–517.
- Sanner, S., and Boutilier, C. 2009. Practical solution techniques for first-order mdps. *Artif. Intell.* 173(5-6):748–788.
- Sanner, S. 2010. Relational Dynamic Influence Diagram Language (RDDL): Language Description.
- Shen, W.; Trevizan, F.; Toyer, S.; Thiébaux, S.; and Xie, L. 2019. Guiding Search with Generalized Policies for Probabilistic Planning. In *Proc. of 12th Annual Symp. on Combinatorial Search (SoCS)*.
- Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.; van den Driessche, G.; Schrittwieser, J.; Antonoglou, I.; Panneershelvam, V.; Lanctot, M.; Dieleman, S.; Grewe, D.; Nham, J.; Kalchbrenner, N.; Sutskever, I.; Lillicrap, T. P.; Leach, M.; Kavukcuoglu, K.; Graepel, T.; and Hassabis, D. 2016. Mastering the game of go with deep neural networks and tree search. *Nature* 529(7587):484–489.
- Sorg, J., and Singh, S. P. 2009. Transfer via soft homomorphisms. In *8th International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS 2009)*, Budapest, Hungary, May 10-15, 2009, Volume 2, 741–748.

Taylor, M. E., and Stone, P. 2009. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research* 10:1633–1685.

Toyer, S.; Trevizan, F. W.; Thiébaux, S.; and Xie, L. 2018. Action schema networks: Generalised policies with deep learning. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, Louisiana, USA, February 2-7, 2018*.

Velickovic, P.; Cucurull, G.; Casanova, A.; Romero, A.; Liò, P.; and Bengio, Y. 2017. Graph attention networks. *CoRR* abs/1710.10903.

Xu, B.; Wang, N.; Chen, T.; and Li, M. 2015. Empirical evaluation of rectified activations in convolutional network. *CoRR* abs/1505.00853.

Younes, H. L. S.; Littman, M. L.; Weissman, D.; and Asmuth, J. 2005. The first probabilistic track of the international planning competition. *J. Artif. Intell. Res.* 24:851–887.

## Appendix

### A Domain Description

We describe the details of the domains presented in the IPPC 2011 and IPPC 2014. The statistics for state fluents ( $\mathcal{F}$ ), non-fluents ( $\mathcal{NF}$ ) and Action ( $\mathcal{A}$ ) for all the domains are show in the Table 5 and Table 6. UP represent  $\mathcal{F}$ ,  $\mathcal{NF}$  and  $\mathcal{A}$  without parameters, Unary represents  $\mathcal{F}$ ,  $\mathcal{NF}$  and  $\mathcal{A}$  with a single parameter and multiple represents  $\mathcal{F}$ ,  $\mathcal{NF}$  and  $\mathcal{A}$  with more than one parameter. Table 7 lists the instance specific number of objects, state variables and action variables for the domain. The domains 1, 2, 3 are used for training, 4 for validation and 5, 6, 7, 8, 9, 10 for testing.

#### Academic Advising

The academic advising domain represents a student at a university trying to complete his/her degree. Some courses are required to be completed to obtain the final degree. Each course is either a basic course or may have prerequisites. The probability of passing a course depends on the number of prerequisites completed (a fixed probability if no prerequisite). The goal is to complete the degree as soon as possible.

#### Crossing Traffic

Crossing Traffic is represents a robot in a grid, with obstacles at a random grid cell at any time. The obstacles (car) start at any cell randomly and move left. The robot aims to plan its path from the starting grid cell to the goal cell while avoiding obstacles.

#### Game of Life

Game of Life domain is represented as a grid where each cell can either be dead or alive. The goal is to keep as many cells alive as possible. The probability of cell death depends on the number of neighbors alive at a particular time, which is non-linear in the number of neighbors alive.

#### Navigation

Navigation represents a robot in a grid world where the aim is to reach a goal cell as quickly as possible. The probability of the robot dying in a particular cell is different, which is specified in the instance file.

#### Skill Teaching

Skill Teaching domain represents a teacher trying to teach a skill to students. Each student has a mastery level in a particular skill. Some skills have pre-conditions, which increase the probability of learning a particular skill. The skill is taught using either hints or multiple-choice questions. The goal is to answer as many questions as possible by the student by learning the required skill.

#### Sysadmin

Sysadmin domain represents computers connected in a network. The probability of a computer shutting down on its own depends on the number of turned-on neighboring computers. The agent can either turn on a computer or leave it as it is. The goal is to maximize the number of computers at a particular time.

Table 5: The statistics related to the domains listing the number of UP (Un-Paramataried), Unary and Multiple Action ( $\mathcal{A}$ ) for each domain.

Domain	UP- $\mathcal{A}$	Unary- $\mathcal{A}$	Multiple- $\mathcal{A}$
Academic Advising	0	1	0
Crossing Traffic	4	0	0
Game of Life	0	0	1
Navigation	4	0	0
Skill Teaching	0	2	0
Sysadmin	0	1	0
Tamarisk	0	2	0
Traffic	0	1	0
Wildfire	0	0	2

#### Tamarisk

Tamarisk domain represents invasive species of plants (Tamarisk) trying to take over native plant species. The plants spread in any direction and try to destroy the native plant species. The agent can either eradicate Tamarisk in a cell or restore the native plant species, each having a different reward. The goal is to minimize the cost of eradication and restoration of the native plant species.

#### Traffic

Traffic domain models the traffic on the road with roads connecting at various intersections. Each road intersection has two traffic light signals combinations of which yield different traffic movement. The agent aims to control the traffic signal (only on the forward sequence) to control the traffic.

#### Wildfire

The wildfire domain represents a forest catching fire. The direction of fire spreading depends on the direction of the wind and also the type of fuel at that point (e.g., grass or wood, etc.). The agent can either choose to put down the fire or cut off the fuel even before the fire happens. The goal is to prevent as many cells as possible, and more reward is provided to protect high priority cells.

### B Variation of $\alpha_{\text{SYMNET}}(0)$ with neighbourhood

To inspect the importance of the neighborhood information in learning a generalized policy for the domains, we perform the study of the neighborhood parameter variation. In the Figure 3, we show the variation of  $\alpha_{\text{SYMNET}}(0)$  with neighbourhood. From the Figure, we observe that message passing for the neighborhood of size 1 yields the best results for most domains, and hence we reported the results with neighborhood 1 in the main paper. In general, we observe that the value of  $\alpha_{\text{SYMNET}}(0)$  first increases and then decreases.

For most instances, the  $\alpha_{\text{SYMNET}}(0)$  is less for neighborhood 0 compared to neighborhood 1, showing that the information regarding the neighbors is necessary for learning a better policy. For example, in domain academic advising,

Table 6: The statistics related to the domains listing the number of UP (Un-Paramataried), Unary and Multiple State Fluents ( $\mathcal{F}$ ) and Non-Fluents ( $\mathcal{NF}$ ) for each domain.

Domain	UP- $\mathcal{F}$	UP- $\mathcal{NF}$	Unary- $\mathcal{F}$	Unary- $\mathcal{NF}$	Multiple- $\mathcal{F}$	Multiple- $\mathcal{NF}$
Academic Advising	0	1	2	5	0	1
Crossing Traffic	0	1	0	4	2	5
Game of Life	0	0	0	0	1	2
Navigation	0	0	0	4	1	6
Skill Teaching	0	0	6	7	0	1
Sysadmin	0	2	1	0	0	1
Tamarisk	0	17	2	0	0	2
Traffic	0	0	3	3	0	3
Wildfire	0	4	0	0	2	2

the neighborhood 1 aggregates information about the pre-requisites for the courses and then prioritizes the courses to take. A similar trend is observed in domain skill teaching, where the information about the pre-condition for the skill plays an important role in learning the skills. For some domains like navigation, neighborhood information is absolutely critical for planning the next move which can be observed from very low values of  $\alpha_{\text{SYMNET}}(0)$  from Figure 3(d). Other domains like wildfire are not affected a lot by neighborhood a lot. This is because the margin between the minimum and maximum rewards is large, and the generalized policy outputs rewards close to the maximum value, which decreases the variation in the value of  $\alpha_{\text{SYMNET}}(0)$ . As we increase the value of neighborhood to 2 and 3, the value of  $\alpha_{\text{SYMNET}}(0)$  tends to fall down for most instances. We hypothesize that the agent overfits to instance-specific policies for the instances it is trained on and hence fails to generalize.

Table 7: The statistics related to the domain instances listing the number of Objects, State Variables and Action Variables for all the instances of the domains. We also give the number of node and edges in the instance graph  $G_{\mathcal{K}}$  (i.e. the graph constructed from the DBN). Domain 1, 2, 3 are used for training, 4 for validation and 5, 6, 7, 8, 9, 10 for testing.

Domain	#Objects	#State Vars	#Action Vars	#Nodes	#Edges	Domain	#Objects	#State Vars	#Action Vars	#Nodes	#Edges
AA 1	10	20	11	10	26	ST 1	2	12	5	2	4
AA 2	10	20	11	10	32	ST 2	2	12	5	2	4
AA 3	15	30	16	15	38	ST 3	4	24	9	4	16
AA 4	15	30	16	15	51	ST 4	4	24	9	4	16
AA 5	20	40	21	20	53	ST 5	6	36	13	6	36
AA 6	20	40	21	20	63	ST 6	6	36	13	6	36
AA 7	25	50	26	25	63	ST 7	7	42	15	7	49
AA 8	25	50	26	25	93	ST 8	7	42	15	7	49
AA 9	30	60	31	30	75	ST 9	8	48	17	8	64
AA 10	30	60	31	30	101	ST 10	8	48	17	8	64
CT 1	9	12	5	9	39	Sys 1	10	10	11	10	24
CT 2	9	12	5	9	39	Sys 2	10	10	11	10	38
CT 3	16	24	5	16	77	Sys 3	20	20	21	20	58
CT 4	16	24	5	16	77	Sys 4	20	20	21	20	77
CT 5	25	40	5	25	127	Sys 5	30	30	31	30	86
CT 6	25	40	5	25	127	Sys 6	30	30	31	30	111
CT 7	36	60	5	36	189	Sys 7	40	40	41	40	118
CT 8	36	60	5	36	189	Sys 8	40	40	41	40	156
CT 9	49	84	5	49	263	Sys 9	50	50	51	50	150
CT 10	49	84	5	49	263	Sys 10	50	50	51	50	196
GOL 1	9	9	10	9	49	Tam 1	12	16	9	8	40
GOL 2	9	9	10	9	49	Tam 2	16	24	9	12	90
GOL 3	9	9	10	9	49	Tam 3	15	20	11	10	52
GOL 4	16	16	17	16	100	Tam 4	20	30	11	15	117
GOL 5	16	16	17	16	100	Tam 5	18	24	13	12	64
GOL 6	16	16	17	16	100	Tam 6	24	36	13	18	144
GOL 7	25	25	26	25	169	Tam 7	21	28	15	14	76
GOL 8	25	25	26	25	169	Tam 8	28	42	15	21	171
GOL 9	25	25	26	25	169	Tam 9	24	32	17	16	88
GOL 10	30	30	31	30	196	Tam 10	32	48	17	24	198
Nav 1	12	12	5	12	55	Tra 1	28	32	5	28	84
Nav 2	15	15	5	15	71	Tra 2	28	32	5	28	84
Nav 3	20	20	5	20	99	Tra 3	40	44	5	40	120
Nav 4	30	30	5	30	155	Tra 4	40	44	5	40	120
Nav 5	30	30	5	30	151	Tra 5	52	56	5	52	156
Nav 6	40	40	5	40	209	Tra 6	52	56	5	52	156
Nav 7	50	50	5	50	267	Tra 7	64	68	5	64	192
Nav 8	60	60	5	60	311	Tra 8	64	68	5	64	192
Nav 9	80	80	5	80	429	Tra 9	76	80	5	76	228
Nav 10	100	100	5	100	547	Tra 10	76	80	5	76	228
Wild 1	9	18	19	9	48	Wild 6	25	50	51	25	156
Wild 2	9	18	19	9	43	Wild 7	30	60	61	30	180
Wild 3	16	32	33	16	87	Wild 8	30	60	61	30	172
Wild 4	16	32	33	16	97	Wild 9	36	72	73	36	238
Wild 5	25	50	51	25	156	Wild 10	36	72	73	36	230

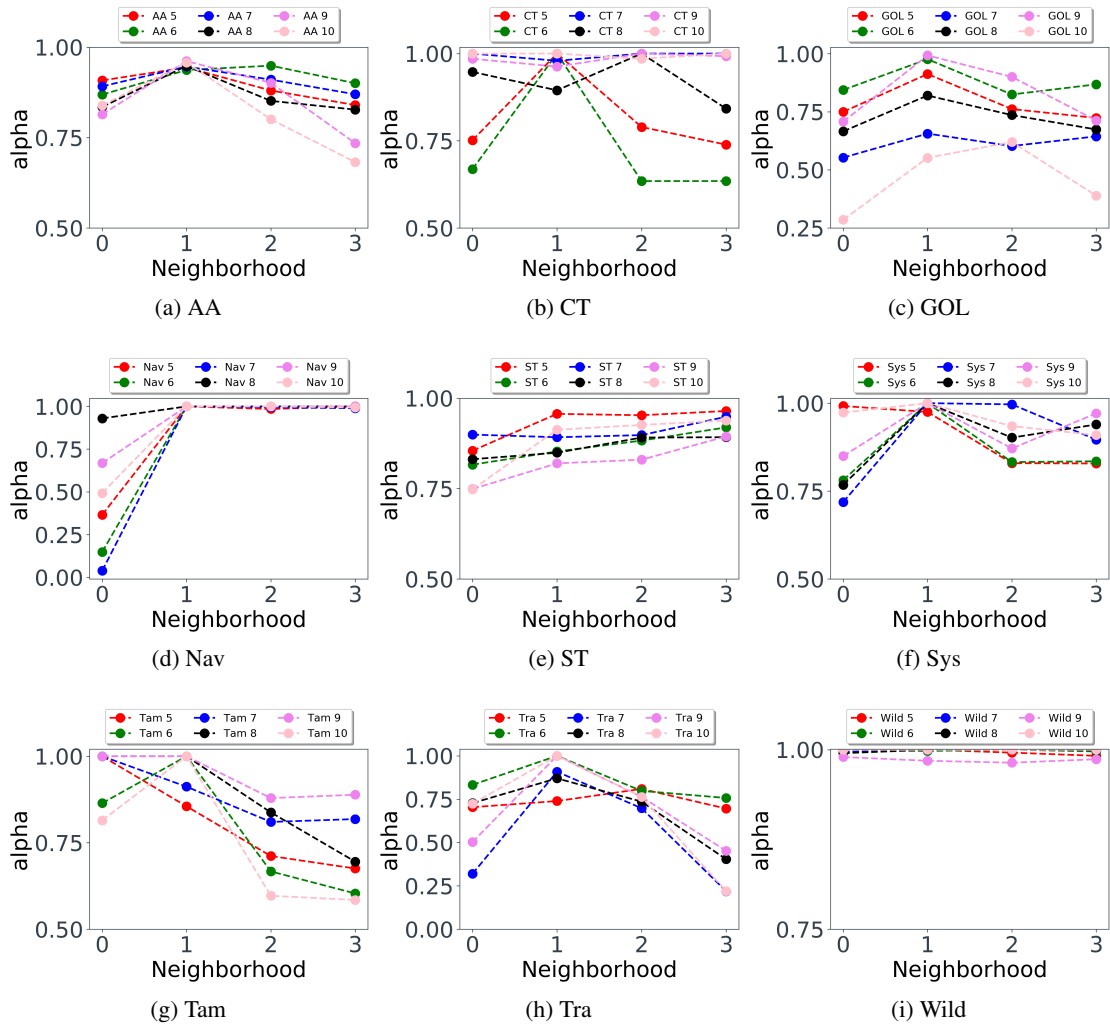


Figure 3: Variation of  $\alpha_{SYMNET}(0)$  with neighbourhood. [Larger is better]