

DeepCubeAF: A Foundation Model for Generalizable Pathfinding Heuristics

Anonymous authors

Paper under double-blind review

Keywords: foundation models, generalization in planning, deep reinforcement learning, pathfinding

Summary

We present DeepCubeAF, a method for training a foundation model for heuristic functions that generalize across pathfinding domains. Our approach uses a domain generator to generate training data, trains a graph neural network (GNN) using deep reinforcement learning and hindsight experience replay, and uses batch weighted A* search to solve problems. Our experiments show that DeepCubeAF consistently solves more problems than foundation models based on supervised learning, is competitive with traditional domain-independent planners, and can solve problems from domains not seen during training.

Contribution(s)

1. We introduce DeepCubeAF, a foundation model trained with reinforcement learning to learn heuristic functions for solving pathfinding problems.
Context: Previous methods rely on supervised learning and, therefore, assume training instances can already be solved. Traditional domain-independent methods do not make use of machine learning, and thus must rely only on hand-designed heuristics.
2. We demonstrate that the resulting model can generalize to previously unseen domains by applying without additional fine-tuning.
Context: This shows that deep reinforcement learning is a viable approach for creating a domain-independent solver. This will allow practitioners that are not familiar with machine learning algorithms to solve pathfinding problems with machine learning models.

DeepCubeAF: A Foundation Model for Generalizable Pathfinding Heuristics

Anonymous authors

Paper under double-blind review

Abstract

1 Pathfinding problems can be found in fields such as robotics, mathematics, chemistry,
 2 and program synthesis, where the objective of pathfinding is to find a sequence of ac-
 3 tions that transforms a given start state into a goal state. Recently, deep reinforcement
 4 learning (DRL) has emerged as a promising method for automatically training domain-
 5 specific heuristic functions to solve these problems in a largely domain-independent
 6 fashion. However, these approaches often require retraining for even a slight change in
 7 domain, resulting in significant resource and time inefficiencies. While existing ap-
 8 proaches use supervised learning to learn generalizable heuristics to handle unseen
 9 domains, they are limited by the need to obtain supervised labels. We draw inspi-
 10 ration from domain randomization in reinforcement learning to handle these limita-
 11 tions and the DeepCubeA algorithm and introduce DeepCubeA for foundation models
 12 (DeepCubeAF). DeepCubeAF trains a heuristic function across randomly generated do-
 13 mains using reinforcement learning and uses this trained heuristic function with batch
 14 weighted A* search to solve problems. Our model consistently shows better generaliz-
 15 ability than the existing foundation model for both seen and unseen domains. This work
 16 represents a step toward training robust, generalizable models and providing access to
 17 these models to experts across various fields.

18 1 Introduction

19 Pathfinding aims to find a sequence of actions that forms a path from a given start state to a given
 20 goal state while minimizing the total path cost. Pathfinding problems, many of which are also stud-
 21 ied in reinforcement learning, are prevalent across computer science, robotics, mathematics, and
 22 the natural sciences. Heuristic search, one of the most prominent approaches to solving pathfinding
 23 problems, relies on a heuristic function which estimates the “cost-to-go,” which is the cost of the
 24 shortest path from a given state to a nearest goal state. Recently, deep reinforcement learning (DRL)
 25 (Sutton & Barto, 2018) based methods have successfully been used to learn domain-specific heuris-
 26 tics in a largely domain-independent fashion to solve pathfinding problems such as puzzle solving
 27 (Agostinelli et al., 2019; 2024), chemical synthesis (Chen et al., 2020), and quantum algorithm com-
 28 pilation (Zhang et al., 2020). However, training such heuristic functions using deep neural networks
 29 (DNNs) (Schmidhuber, 2015) can take days and requires substantial computational resources, of-
 30 ten necessitating retraining for even minor domain changes. These hardware and time constraints
 31 limit accessibility of these methods for researchers not familiar with machine learning algorithms
 32 and hardware and present a bottleneck when attempting to solve problems across many different
 33 domains.

34 These challenges with hardware and time constraints when using DRL to solve pathfinding prob-
 35 lems are common across applications of machine learning, such as those encountered in computer
 36 vision and natural language processing. To address these challenges, foundation models have been
 37 developed, which are large DNNs trained on large, diverse, datasets that can then be adapted to new

tasks with minimal to no fine-tuning (Szegedy et al., 2015; Devlin et al., 2019; Radford et al., 2019). Based on the success of foundation models for other applications of machine learning, we will create a foundation model for heuristic functions for pathfinding problems.

To create our foundation model, we will use the planning domain description language (PDDL) (McDermott, 2000), which is a description language commonly used by the planning community to represent pathfinding problems. A PDDL description of the domain, along with a PDDL description of the problem instance, which consists of a start state and a goal, will be given to a graph neural network (GNN), which represents the heuristic function. To generate training data, we will build on PDDL FUSE (Khandelwal et al., 2024) to generate domains and hindsight experience replay (HER) (Andrychowicz et al., 2017) to generate training problem instances. To train the heuristic function, we will build on the DeepCubeA algorithm (Agostinelli et al., 2019) by using approximate value iteration (AVI) (Bertsekas & Tsitsiklis, 1996). To solve problem instances in a given domain, we will also build on DeepCubeA by using batch weighted A* search (Hart et al., 1968; Pohl, 1970) with the trained heuristic function. We call this algorithm DeepCubeA for foundation models (DeepCubeAF). Our experiments show that DeepCubeAF can generalize to domains not seen during training and outperforms foundation models trained only using supervised learning.

2 Related Work

2.1 Generalization in Pathfinding Problems

Domain-independent planners, such as the fast downward planner (Helmert, 2006), rely on the automated construction of heuristic functions based on a PDDL description of a domain. These methods include the fast forward (Hoffmann & Nebel, 2001), causal graph, and goal count heuristics. However, experiments have shown that these methods do not perform as well as learned heuristic functions (Agostinelli et al., 2024). Large language models (LLMs), have been used to solve pathfinding via in-context learning (Sermanet et al., 2023; Li et al., 2023; Silver et al., 2023); however, they face notable challenges that include syntax errors and lack inherent search capabilities. Although, future work using chain-of-thought reasoning may be able to overcome this.

Chen et al. (2024) introduced three novel graph representations for planning tasks using GNNs to learn domain-independent heuristics. Their approach mitigates issues with large grounded GNNs by leveraging lifted representations and demonstrates superior generalization to larger problems compared to models like STRIPS-HGN (a hypergraph network model that learns domain-independent planning heuristics directly from the delete-relaxed representation of STRIPS problems using a recurrent encode-process-decode architecture) (Shen et al., 2020). However, it faces scalability issues with large graph construction. Building on this, Chen et al. (2024) proposed the GOOSE framework using GNNs with novel grounded and lifted graph representations for classical planning to learn a generalized heuristic function with supervised learning. Their heuristics outperform STRIPS-HGN and hFF (the Fast Forward heuristic (Hoffmann & Nebel, 2001), which computes cost-to-go estimates by generating relaxed plans that ignore delete effects) in various domains, however, since they use supervised learning, their approach assumes that ability to find optimal solutions for all training instances. However, this is not possible for many pathfinding tasks without significant domain-dependent effort (Agostinelli et al., 2024). Additionally, (Toyer et al., 2018) utilized GNNs to improve coverage and plan quality in classical planning tasks through new graph representations, though their approach only generalizes across a subset of test domains.

2.2 GOOSE

Our work also compares with an existing generalizable method for solving planning domains. (Chen et al., 2024) uses the STRIPS Learning Graph (SLG) representation, which encodes states, actions, and their relationships into a structured graph. SLG captures domain knowledge by representing preconditions and effects and deleting effects of actions as graph edges, enabling the model to learn heuristics through Graph Neural Networks (GNNs).

86 This approach combines domain-specific training with domain-independent methods, leveraging the
 87 SLG’s ability to generalize within a domain. The GNN propagates information through the SLG to
 88 compute heuristic values, allowing the model to handle larger instances of the same domain without
 89 retraining.

90 Although primarily domain-specific, the SLG-based model aligns with foundational principles by
 91 utilizing structural representations and scalable learning techniques. It demonstrates how graph-
 92 based representations can improve adaptability in classical planning tasks, particularly for domains
 93 like n-puzzle.

94 2.3 Domain Randomization in Reinforcement Learning

95 Domain randomization has proven highly effective in reinforcement learning (RL) for training
 96 agents that generalize to unseen environments. By training on multiple randomized variants of a
 97 domain, RL agents learn to handle discrepancies in states, transitions, or dynamics. For instance,
 98 [Mehta et al. \(2020\)](#) introduce Active Domain Randomization, which adaptively selects challenging
 99 environments for training, while [Ajani et al. \(2023\)](#) demonstrate improved robustness by randomiz-
 100 ing physical parameters such as friction in robotics. These successes highlight the power of exposing
 101 learning systems to diverse environments, a principle that can be translated to automated planning
 102 domains.

103 3 Preliminaries

104 3.1 Foundation Models

105 Foundation models are deep learning models pre-trained on extensive, diverse datasets. Once
 106 trained, F_θ is expected to generalize across different domains and tasks with or without fine-tuning.
 107 Foundation models have significant applications in computer vision [Szegedy et al. \(2015\)](#) and natu-
 108 ral language processing (NLP) [Devlin et al. \(2019\)](#); [Radford et al. \(2019\)](#).

109 **Definition:** A foundation model F_θ is trained to minimize the loss function L over a dataset D :

$$\theta^* = \arg \min_{\theta} L(F_\theta(x), x) \quad \text{for } x \in D$$

110 3.2 Pathfinding

111 A pathfinding domain can be represented by a weighted directed graphs where nodes represent
 112 states, edges represent transitions between states (actions), and weights represent transition costs.
 113 Given a domain, a pathfinding problem instance is defined by a start state and a set of goal states.
 114 A solution to a pathfinding problem is a sequence of actions that forms a path between the start
 115 state and a goal state, with preference for solutions with cheaper path costs, where the path cost is
 116 the sum of transition costs along a given path. Traditionally, pathfinding employs heuristic search
 117 algorithms, which make use of a heuristic function that maps states to an estimate of the cost-to-go.

118 **Approximate Value Iteration** Given a tabular representation for the heuristic function, value it-
 119 eration can be used to find h^* , which maps states to the cost of a shortest path using the following
 120 equation as an update rule:

$$h'(s) = \min_{a \in \mathcal{A}} (c(s, a) + h(T(s, a))) \quad (1)$$

121 where $h(s)$ represents the current estimate of the cost-to-go, $h'(s)$ is the updated estimate, $c(s, a)$
 122 is the cost of taking action a in state s , and $T(s, a)$ is the resulting state after applying action a .
 123 Since pathfinding problems are deterministic, value iteration can use a deterministic state transition
 124 function instead of transition probabilities.

125 Tabular methods become impractical for domains with large state spaces, which are common in real-
 126 world problems. To address this, we use approximate value iteration (AVI), where a parameterized
 127 function approximates value iteration updates. In AVI, a DNN with parameters θ is used to represent
 128 the heuristic function. The network refines its estimates for the cost-to-go values by minimizing the
 129 loss:

$$L(\theta) = \left(\min_{a \in \mathcal{A}} (c(s, a) + h_{\theta^-}(T(s, a))) - h_{\theta}(s) \right)^2 \quad (2)$$

130 where $h_{\theta}(s)$ approximates the cost-to-go and θ^- denotes the parameters of a target network (Mnih
 131 et al., 2015), which is periodically updated to θ .

132 3.3 Batch Weighted A* Search

133 A* search (Hart et al., 1968) is a widely used algorithm for solving pathfinding problems. A* search
 134 maintains a search tree, where nodes represent states and edges represent transitions between states.
 135 Nodes are expanded according to a priority, f , shown in Equation 3:

$$f(n) = g(n) + h(n) \quad (3)$$

136 where $g(n)$ is the path cost from the starting node to node n and $h(n)$ is the cost-to-go estimate
 137 from the state associated with node n to a closest goal state provided by a heuristic function. The
 138 algorithm terminates when a node associated with a goal state is selected for expansion.

139 In order to trade potentially faster search times with potentially longer path costs a weight, λ , can
 140 be used to reduce the contribution of, g , to the total cost, as shown in Equation 4, where $0 \leq \lambda \leq 1$.
 141 Furthermore, when the heuristic function is represented by a DNN, parallelism provided by graphics
 142 processing units (GPUs) can be exploited by expanding N nodes at a time instead of just one using
 143 batch weighted A* search (BWAS) (Agostinelli et al., 2019; Li et al., 2022).

$$f(n) = \lambda g(n) + h(n) \quad (4)$$

144 3.4 Planning Domain Definition Language (PDDL)

145 The Planning Domain Definition Language (PDDL) provides a formal framework for representing
 146 planning problems, where an agent transitions through states by executing actions to achieve a goal.
 147 A PDDL problem consists of a *domain* \mathcal{D} and a *problem* \mathcal{P} . The domain \mathcal{D} defines the set of possible
 148 actions and how they effect a given state, while the problem \mathcal{P} specifies an initial state s_0 and a goal
 149 condition G .

150 A PDDL domain is defined by:

- 151 • **Objects** \mathcal{O} : The entities in the domain.
- 152 • **Predicates** \mathcal{P} : Boolean functions describing state conditions.
- 153 • **Actions** \mathcal{A} : Each action $a \in \mathcal{A}$ is defined by:
 - 154 – **Parameters**: A set of objects required for execution.
 - 155 – **Preconditions** $\text{Pre}(a) \subseteq \mathcal{P}$: Conditions that must hold before execution.
 - 156 – **Effects** $\text{Eff}(a) \subseteq \mathcal{P}$: State changes resulting from execution.

157 A problem instance \mathcal{P} consists of an initial state $s_0 \subseteq \mathcal{P}$ and a goal condition $G \subseteq \mathcal{P}$. The objective
 158 is to find a sequence of actions $\pi = (a_1, a_2, \dots, a_n)$ such that executing π from s_0 results in a state
 159 satisfying G .

160 3.5 PDDL FUSE: Generating Diverse Planning Domains

161 Taking inspiration from domain randomization techniques in reinforcement learning (Mehta et al.,
 162 2020; Ajani et al., 2023), PDDL FUSE (Khandelwal et al., 2024) generates planning domains in

163 PDDL format by randomly modifying and combining (fusing) existing domains. PDDL_{FUSE} does
 164 this by by: 1) merging predicate sets and action schemas from multiple domains; 2) randomly adding
 165 or removing preconditions and/or effects; 3) introducing controlled negation and reversibility flags
 166 for selected predicates. Empirical results show that standard domain-independent planners often
 167 struggle to solve the more challenging PDDL_{FUSE} domains.

168 4 Methodology

169 4.1 Enhanced Domain and Problem Instance Generation

170 **Expanded Base Domains and Problem Files.** Our modifications to PDDL_{FUSE} begin by signif-
 171 icantly increasing the set of base domains and problem files. This expansion ensures more diverse
 172 set of initial problem files per generated domains. When fusing domains, we preserve predicate and
 173 action names by applying systematic renaming, ensuring there is no overlap among symbols across
 174 different domains. The final fused domain thus contains a merged set of predicates and actions with
 175 randomized modifications to preconditions and effects.

176 **Adaptive Negation.** A key step in domain fusion involves adding or removing preconditions and
 177 effects from actions. We introduce an adaptive negation mechanism that selectively identifies pred-
 178 icates eligible for negation in the newly fused domain’s problem file. Specifically, a predicate is
 179 deemed eligible if it is not present in the critical preconditions of actions—since negating such pred-
 180 icates could prevent actions from being executed and render the problem unsolvable. By applying
 181 negation with a user-controlled probability only to these non-essential predicates, our approach pre-
 182 serves the core functionality of the domain while still introducing useful variability. By carefully
 183 targeting the negated predicates, we improve the consistency of generated actions and the solvability
 184 of resultant problems.

185 4.2 Revised Action Generation Technique

186 Previously, PDDL_{FUSE} created problem files by randomly selecting actions from the fused domain
 187 and subsequently pruning to form a valid plan. This random approach was often inefficient, requiring
 188 substantial computation to verify solvability and identify relevant actions. We replace it with a more
 189 planner-aligned technique designed to efficiently produce action sequences and solvable planning
 190 problems.

191 **Planner-Aligned Sampling.** Instead of choosing actions uniformly at random, we adopt a tech-
 192 nique akin to a forward search used in classical planners. We begin with an initial state derived from
 193 the domain objects and predicates, then dynamically select only the set of applicable actions—those
 194 whose preconditions are satisfied in the current state. Among these applicable actions, we apply a
 195 stochastic selection (e.g., uniform or heuristic-based) to pick one action. We then apply this action
 196 to reach a successor state.

197 **Controlled Sequence Generation.** We repeat the above process for a user-configurable number
 198 of steps (or until no applicable actions remain). This controlled, state-based generation ensures that
 199 we continuously build a coherent action sequence, thereby avoiding many dead-end paths. Once the
 200 final state is reached, we designate a subset of its predicates as the *goal* conditions. The resulting
 201 domain and problem pair therefore includes a sequence of actions known to transition the initial
 202 state to the goal state. This process is summarized in Algorithm 1. This approach is similar to that
 203 of hindsight experience replay (HER) [Andrychowicz et al. \(2017\)](#) in that the state at the end of a
 204 sequence of actions is used to create a goal.

205 **Computational Benefits.** Aligning the action generator with planner logic significantly lowers the
 206 computational overhead. Instead of repeatedly checking for plan validity after random selections,
 207 the state-based approach validates each action as it is chosen, mitigating expensive re-planning or

Algorithm 1 Revised Action Generator**Require:** Fused domain \mathcal{D} , maxSteps M

```

1: state  $\leftarrow$  randomlyInstantiate( $\mathcal{D}$ )
2: for  $i = 1$  to  $M$  do
3:   applicableActions  $\leftarrow \{a \in \mathcal{A} \mid \text{pre}(a) \subseteq \text{initializeState}\}$ 
4:   if applicableActions =  $\emptyset$  then
5:     break ▷ No more valid actions
6:   end if
7:    $a^* \leftarrow$  chooseAction(applicableActions)
8:   state  $\leftarrow$  apply( $a^*$ , state)
9: end for
10: goalPreds  $\leftarrow$  selectSubset(state)
11: return ( $\mathcal{D}$ , initialState, goalPreds)

```

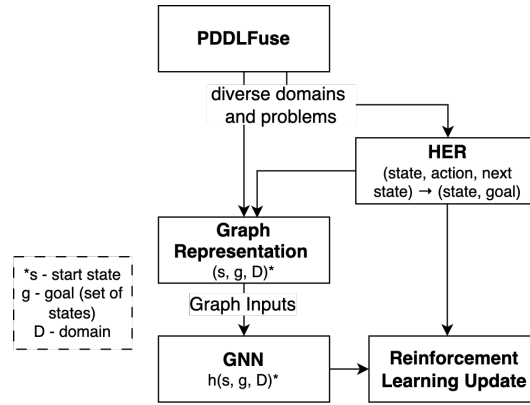


Figure 1: Overview of the training pipeline. PDDLFuse generates diverse domains and problem instances from which states and goals (which is a set of states) are derived. HER is used to generate problem instances, and a graph-based representation encodes each state, goal (a set of states), and domain as a graph. The GNN maps this representation to an estimate of the cost-to-go.

208 pruning phases. This streamlined method is especially beneficial when large numbers of domains or
 209 problem files are generated, as it accelerates the overall workflow.

210 4.3 Training the Foundation Model

211 Figure 1 illustrates our training algorithm which combines PDDLFuse, HER, GNNs, and rein-
 212 forcement learning.

213 **State and Goal Generation.** We use PDDLFuse to create diverse planning domains \mathcal{D} and cor-
 214 responding problem instance Π . Each state, s , is associated with a domain and problem instance.
 215 We then sample actions under a given exploration policy to produce transitions (s, a, s') . After each
 216 transition, drawing from HER, we create additional goal states, g' , by treating s' (or subsets of its
 217 predicates) as potential new goals. This process increases the amount of training data by using paths
 218 that did not reach the original goal, g .

219 **Graph Representation.** Given a state, goal, and domain tuple (s, g, \mathcal{D}) , we construct a STRIPS
 220 Learning Graph (SLG) to encode the planning task (Chen et al., 2024). Each node corresponds to
 221 grounded predicates from s , with a feature vector indicating whether the predicate is true in the ini-
 222 tial state and whether it is required by the goal g . Edges capture semantic relationships defined by
 223 actions: an edge labeled "pre" links an action to a predicate in its preconditions, while edges labeled

"add" or "del" link the action to its positive or negative effects. This concise, structured representation encapsulates the core dependencies between predicates and actions, providing a rich input for our Graph Neural Network to learn effective heuristic functions. This representation produces a set of inputs for a GNN.

GNN Heuristic Function. A GNN, h_θ , takes (s, g, \mathcal{D}) as input and outputs an estimate of the cost-to-go, $h_\theta(s, g, \mathcal{D})$. During training, we minimize the value-iteration-based loss shown in Equation 5:

$$L(\theta) = \left(\min_{a \in \mathcal{A}} (c(s, a) + h_\theta(T(s, a), g, \mathcal{D})) - h_\theta(s, g, \mathcal{D}) \right)^2 \quad (5)$$

5 Experiments

We use eight canonical planning domains: Blocks, Ferry, Gripper, N-Puzzle, Sokoban, Spanner, VisitAll, and VisitSome. Each domain has distinct predicates and actions. These domains are used during training for GOOSE and with PDDL FUSE to train DeepCubeAF. A ninth domain, Folding, is used as an additional test domain that was not used during training. We also train a domain-specific heuristic function using DeepCubeA and HER on the folding domain to compare how a heuristic function trained for a specific domain compares to that of DeepCubeAF.

We train DeepCubeAF for one million iterations on domains generated by PDDL FUSE using the Adam optimizer with an initial learning rate of 0.001, with decay of 0.9999993, and batch size of 1000. We encode planning states using a STRIPS Learning Graph (SLG) representation (Chen et al., 2024). This graph captures domain predicates, objects, and actions in a structured manner. All learned models (Vanilla DeepCubeA, DeepCubeAF, and GOOSE) use a message-passing neural network with 16 message-passing layers. This setup follows the GOOSE architecture to maintain consistency among learned models.

We compare to GOOSE, which trains a GNN to represent a heuristic function with supervised learning on the eight aforementioned domains. Since DeepCubeAF is based on reinforcement learning, it does not assume the ability to solve the given problem instances. Therefore, DeepCubeAF can train on many more domains and problem instances. We also compare to the fast downward planner with the fast-forward (FF) heuristic (Helmert, 2006), which is a domain-independent planner that does not use machine learning.

To solve problem instances with the trained heuristic function from DeepCubeAF, we use batch weighted A* search (BWAS) with 100 batch size and 0.8 weight. For GOOSE, we report results available in their paper (Chen et al., 2024) as well as when using the GOOSE heuristic function with our implementation BWAS. Finally, to further test the ability of DeepCubeAF to generalize, we do a leave-one-out experiment, where one of the eight training domains is not used in PDDL FUSE and a heuristic function is trained with data generated from PDDL FUSE with the remaining seven.

6 Results

We evaluate each solver on eight domains plus an additional domain, folding, which is not used for training DeepCubeAF and GOOSE. We use the test instances from the (Chen et al., 2024), as shown in Table 1, which also indicates the parameter ranges for each domain. Table 2 shows results for all solvers. DeepCubeAF achieves higher coverage than GOOSE in all domains, including the leave-one-out (LOO) variants of DeepCubeAF that exclude the test domain. GOOSE_ours, which combines the original GOOSE heuristic function with our implementation of BWAS, solves more instances than GOOSE reported by Chen et al. (2024) (which may be due to the batched node expansions) but still solves fewer problem instances than DeepCubeAF and DeepCubeAF(LOO). DeepCubeAF outperforms DeepCubeAF(LOO), which suggests that additional domain diversity in

training contributes to better generalization. DeepCubeAF solves more problem instances than the fast downward planner in the majority of cases. In the Spanner domain, DeepCubeAF solves 18% of the instances, DeepCubeAF(LOO) solves 7%, whereas FD(FF) and the GOOSE models do not solve any problems.

Neither DeepCubeAF nor GOOSE uses the folding domain during training. DeepCubeA, trained only on folding, attains the highest coverage on folding, as shown in Table 3. Among the domain-independent approaches, DeepCubeAF solves 60% of problem instances while GOOSE and the fast downward planner solve 38%.

domain	testing
blocks (90)	$b \in [15, 100]$
ferry (90)	$l, c \in [15, 100]$
gripper (18)	$b \in [15, 100]$
n-puzzle (50)	$n \in [5, 9]$
sokoban (90)	$n \in [8, 12]$
spanner (90)	$s, n \in [15, 100]$
visitall (90)	$n \in [15, 100]$
visitsome (90)	$n \in [15, 100]$
domain-folding (90)	length, folds $\in [15, 100]$

Table 1: Domains, parameter ranges, and number of instances for testing.

7 Discussion

DeepCubeAF uses reinforcement learning to train its heuristic function while GOOSE uses supervised learning. Though supervised learning provides exact labels, reinforcement learning does not assume that the problems it sees can already be solved. In fact, many problems generated by PDDL-FUSE cannot be solved with domain-independent planners (Khandelwal et al., 2024), which means they would not be able to be used in a supervised learning setting. However, reinforcement learning will most likely need more examples and longer training times compared to supervised learning since reinforcement learning iteratively updates its own labels. We believe DeepCubeAF was able to outperform GOOSE due to the fact that more diverse training data can be seen during training because it does not assume that labels are given.

One benefit of domain-independent planners that are not based on learning, like the fast-downward planner, is that heuristics can quickly be computed from PDDL files. However, experiments on the folding domain, which was not seen during training for DeepCubeAF, shows that foundation models based on deep reinforcement learning can outperform traditional domain-independent planners. Furthermore, our leave-one-out experiments show DeepCubeAF outperforming the fast-downward planner on three domains. Future research is needed to better define the set of domains for which we can expect foundation models to outperform traditional domain-independent planners.

Future work will investigate how to improve DeepCubeAF, especially for domains not used during training. Foundation models for NLP have showed that larger models and more diverse training data results in better performance. Therefore, one improvement can come from improving the diversity of domains generated by PDDL-FUSE. This can possibly be done by leveraging LLMs to generate new PDDL domains instead of just combining existing domains. Another improvement can come from simply training larger models. Finally, given a test domain, limited finetuning can be done. Future work will investigate the effect of finetuning on performance.

Domain	Solver	Len	Nodes	Secs	Nodes/Sec	Solved
blocks	FD(FF)	366.9	4.74E+05	7.7	1.45E+04	74%
	GOOSE	-	-	-	-	10%
	GOOSE_ours	168.25	1.19E+06	65.73	1.81E+04	18%
	DCAF(LOO)	249.15	4.45E+06	228.21	1.95E+04	44%
	DCAF	405.12	1.89E+07	492.14	3.84E+04	80%
ferry	FD(FF)	245.47	4.78E+05	2.01	1.29E+04	95.5%
	GOOSE	-	-	-	-	31%
	GOOSE_ours	148.75	1.96E+06	86.96	2.27E+04	40%
	DCAF(LOO)	175.09	2.32E+06	148.14	1.56E+04	58%
	DCAF	239.21	4.42E+06	268.68	1.64E+04	92%
gripper	FD(FF)	239	1.67E+04	0.06	6.88E+04	100%
	GOOSE	-	-	-	-	28%
	GOOSE_ours	173.13	9.18E+05	58	1.58E+04	44%
	DCAF(LOO)	175.77	1.32E+06	89.2	1.48E+04	50%
	DCAF	220.1	1.51E+06	142.2	1E+04	89%
n-puzzle	FD(FF)	1529.54	1.59E+06	61.84	1.55E+04	70%
	GOOSE	-	-	-	-	12%
	GOOSE_ours	798.68	1.84E+07	110.4	1.66E+05	38%
	DCAF(LOO)	1013.88	8.62E+06	82.4	1E+04	50%
	DCAF	1414.23	2.06E+07	127.63	1.61E+05	72%
sokoban	FD(FF)	100.38	5.11E+05	7.61	3.22E+04	98%
	GOOSE	-	-	-	-	50%
	GOOSE_ours	104.59	1.8E+06	116.9	1.53E+04	63%
	DCAF(LOO)	103.4	2.54E+06	163.8	1.5E+04	65.5%
	DCAF	100.72	5.17E+06	258.4	2E+04	97%
spanner	FD(FF)	-	-	-	-	-
	GOOSE	-	-	-	-	-
	GOOSE_ours	-	-	-	-	-
	DCAF(LOO)	790.3	3.2E+06	210.0	1.52E+04	7%
	DCAF	765.9	1.32E+07	490.5	2.6E+04	18%
visitall	FD(FF)	1842.14	2.77E+07	252.37	2.89E+04	8%
	GOOSE	-	-	-	-	18%
	GOOSE_ours	1950.32	3.20E+07	340.02	9.41E+04	20%
	DCAF(LOO)	1800.91	4.5E+07	480	9.37E+04	27%
	DCAF	1850.07	5.3E+07	590.4	8.97E+04	45%
visitsome	FD(FF)	2110.69	1E+06	30	2.89E+04	29%
	GOOSE	-	-	-	-	81%
	GOOSE_ours	2100.45	9.00E+06	300.0	3.00E+04	64%
	DCAF(LOO)	2000.72	1.10E+07	330.0	3.33E+04	80%
	DCAF	1950.30	2.00E+07	550.0	3.64E+04	85%

Table 2: Comparison of DeepCubeAF (DCAF) with the domain-independent Fast Downward planner (FF heuristic), the GOOSE GNN-based heuristic learner, and the single-domain DCAF(LOO) is a leave-one-out experiment that does not use the given domain during training. Metrics include solution length, node expansions, search time (seconds), nodes expanded per second, and the percentage of solved problems. Note: All other metrics are computed based on solved instances only.

Domain	Solver	Len	Nodes	Secs	Nodes/Sec	Solved
folding	FD(FF)	512.1	8.40E+05	58.47	1.44E+04	38%
	DeepCubeA	420.5	9.00E+06	89.56	1.01E+05	90%
	GOOSE_ours	486.2	2.10E+06	120.20	1.75E+04	38%
	DCAF	600.8	4.80E+06	297.72	1.61E+04	60%

Table 3: Comparison of DeepCubeAF (DCAF) with the domain-independent Fast Downward planner (FF heuristic), the GOOSE GNN-based heuristic learner, and DeepCubeA, which is trained only on the folding domain. The folding domain is not seen during training for DCAF and GOOSE. Metrics include solution length, node expansions, search time (seconds), nodes expanded per second, and the percentage of solved problems. Note: All other metrics are computed based on solved instances only.

8 Conclusion

This paper introduces DeepCubeAF, the first foundation model based on deep reinforcement learning for learning a heuristic function that generalizes across pathfinding domains. DeepCubeAF is trained on randomly generated domains and problem instances and demonstrates improved generalization compared to heuristic functions trained with supervised learning and competitive performance with traditional domain-independent heuristic planners. Experimental results show that DeepCubeAF generalizes to both domains used during training and unseen domains.

References

- Forest Agostinelli, Stephen McAleer, Alexander Shmakov, and Pierre Baldi. Solving the Rubik’s cube with deep reinforcement learning and search. *Nature Machine Intelligence*, 1(8):356–363, 2019.
- Forest Agostinelli, Rojina Panta, and Vedant Khandelwal. Specifying goals to deep neural networks with answer set programming. In *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 34, pp. 2–10, 2024.
- Oladayo S Ajani, Sung-ho Hur, and Rammohan Mallipeddi. Evaluating domain randomization in deep reinforcement learning locomotion tasks. *Mathematics*, 11(23):4744, 2023.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In *Advances in Neural Information Processing Systems*, pp. 5048–5058, 2017.
- Dimitri P Bertsekas and John N Tsitsiklis. *Neuro-dynamic programming*. Athena Scientific, 1996. ISBN 1-886529-10-8.
- Binghong Chen, Chengtao Li, Hanjun Dai, and Le Song. Retro*: learning retrosynthetic planning with neural guided A* search. In *International Conference on Machine Learning*, pp. 1608–1616. PMLR, 2020.
- Dillon Z Chen, Sylvie Thiébaux, and Felipe Trevizan. Learning domain-independent heuristics for grounded and lifted planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 20078–20086, 2024.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.

- 330 Peter E Hart, Nils J Nilsson, and Bertram Raphael. A formal basis for the heuristic determination
331 of minimum cost paths. *IEEE transactions on Systems Science and Cybernetics*, 4(2):100–107,
332 1968.
- 333 Malte Helmert. The fast downward planning system. *Journal of Artificial Intelligence Research*, 26:
334 191–246, 2006.
- 335 Jörg Hoffmann and Bernhard Nebel. The ff planning system: Fast plan generation through heuristic
336 search. *Journal of Artificial Intelligence Research*, 14:253–302, 2001.
- 337 Vedant Khandelwal, Amit Sheth, and Forest Agostinelli. Pddl-fuse: A tool for generating diverse
338 planning domains. *arXiv preprint*, 2024.
- 339 Tianhua Li, Ruimin Chen, Borislav Mavrin, Nathan R Sturtevant, Doron Nadav, and Ariel Felner.
340 Optimal search with neural networks: Challenges and approaches. In *Proceedings of the Interna-*
341 *tional Symposium on Combinatorial Search*, volume 15, pp. 109–117, 2022.
- 342 Yuliang Li, Nitin Kamra, Ruta Desai, and Alon Halevy. Human-centered planning. *arXiv preprint*
343 *arXiv:2311.04403*, 2023.
- 344 Drew M McDermott. The 1998 ai planning systems competition. *AI magazine*, 21(2):35–35, 2000.
- 345 Bhairav Mehta, Manfred Diaz, Florian Golemo, Christopher J Pal, and Liam Paull. Active domain
346 randomization. In *Conference on Robot Learning*, pp. 1162–1176. PMLR, 2020.
- 347 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle-
348 mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level
349 control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 350 Ira Pohl. Heuristic search viewed as path finding in a graph. *Artificial intelligence*, 1(3-4):193–204,
351 1970.
- 352 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
353 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 354 Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural networks*, 61:85–117,
355 2015.
- 356 Pierre Sermanet, Tianli Ding, Jeffrey Zhao, Fei Xia, Debidatta Dwibedi, Keerthana Gopalakrish-
357 nan, Christine Chan, Gabriel Dulac-Arnold, Sharath Maddineni, Nikhil J Joshi, et al. Robovqa:
358 Multimodal long-horizon reasoning for robotics. *arXiv preprint arXiv:2311.00899*, 2023.
- 359 William Shen, Felipe Trevizan, and Sylvie Thiébaux. Learning domain-independent planning
360 heuristics with hypergraph networks. In *Proceedings of the International Conference on Auto-*
361 *mated Planning and Scheduling*, volume 30, pp. 574–584, 2020.
- 362 Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B Tenenbaum, Leslie Pack Kaelbling, and
363 Michael Katz. Generalized planning in pddl domains with pretrained large language models.
364 *arXiv preprint arXiv:2305.11014*, 2023.
- 365 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 366 Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Du-
367 mitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In
368 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- 369 Sam Toyer, Felipe Trevizan, Sylvie Thiébaux, and Lexing Xie. Action schema networks: Gener-
370 alised policies with deep learning. In *Proceedings of the AAAI Conference on Artificial Intelli-*
371 *gence*, volume 32, 2018.
- 372 Yuan-Hang Zhang, Pei-Lin Zheng, Yi Zhang, and Dong-Ling Deng. Topological quantum compiling
373 with reinforcement learning. *Physical Review Letters*, 125(17):170501, 2020.