
Learning Graph Search Heuristics

Anonymous Author(s)

Anonymous Affiliation

Anonymous Email

Abstract

1 Searching for a path between two nodes in a graph is one of the most well-studied
2 and fundamental problems in computer science. In numerous domains such as
3 robotics, AI, or biology, practitioners develop search heuristics to accelerate their
4 pathfinding algorithms. However, it is a laborious and complex process to hand-
5 design heuristics based on the problem and the structure of a given use case. Here
6 we present PHIL (Path Heuristic with Imitation Learning), a novel neural archite-
7 cture and a training algorithm for discovering graph search and navigation heuristics
8 from data by leveraging recent advances in imitation learning and graph represen-
9 tation learning. At training time, we aggregate datasets of search trajectories and
10 ground-truth shortest path distances, which we use to train a specialized graph
11 neural network-based heuristic function using backpropagation through steps of
12 the pathfinding process. Our heuristic function learns graph embeddings useful
13 for inferring node distances, runs in constant time independent of graph sizes, and
14 can be easily incorporated in an algorithm such as A* at test time. Experiments
15 show that PHIL reduces the number of explored nodes compared to state-of-the-art
16 methods on benchmark datasets by 58.5% on average, can be directly applied in
17 diverse graphs ranging from biological networks to road networks, and allows for
18 fast planning in time-critical robotics domains.
19

20 1 Introduction

21 Search heuristics are essential in several domains, including robotics, AI, biology, and chemistry [1–
22 6]. For example, in robotics, complex robot geometries often yield slow collision checks, and search
23 algorithms are constrained by the robot’s onboard computation resources, requiring well-performing
24 search heuristics that visit as few nodes as possible [1, 4]. In AI, domain-specific search heuristics
25 are useful for improving the performance of inference engines operating on knowledge bases [3, 5].
26 Search heuristics have been previously also developed to reduce search efforts in protein-protein
27 interaction networks [6] and in planning chemical reactions that can synthesize target chemical
28 products [2]. This broad set of applications underlines the importance of good search heuristics that
29 are applicable to a wide range of problems.

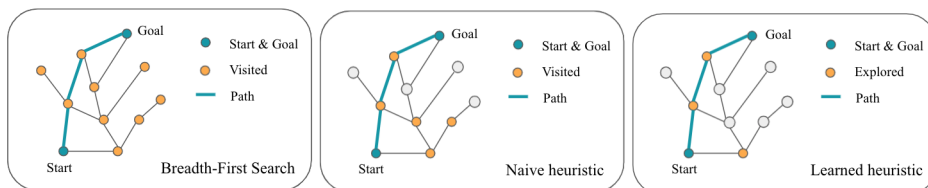


Figure 1: The goal is to navigate (find a path) from the start to the goal node. While BFS visits many nodes to find a start-to-goal path (left), one can use a heuristic based on the features of the nodes (e.g., Euclidean distance) on the graph to reduce the search effort (middle). We propose PHIL to learn a tailored search heuristic for a given graph, capable of reducing the number of visited nodes even further by exploiting the inductive biases of the graph (right).

30 The search task can be formulated as a pathfinding problem on a graph, where given a graph, the
 31 task is to navigate and find a short feasible path from a start node to a goal node, while in the process
 32 visiting as few nodes as possible (Figure 1). The most straightforward approach would be to launch a
 33 search algorithm such as breadth-first search (BFS) and iteratively expand the graph from the start
 34 node until it reaches the goal node. Since BFS does not harness any prior knowledge about the graph,
 35 it usually visits many nodes before reaching the goal, which is expensive in cases such as robotics
 36 where visiting nodes is costly. To visit fewer nodes during the search, one may use domain-specific
 37 information about the graph via a *heuristic* function [7], which allows one to define a distance metric
 38 on graph nodes to prune directions that seem less promising to explore. Unfortunately, coming up
 39 with good search heuristics requires significant domain expertise and manual effort.

40 While there has been significant progress in designing search heuristics, it remains a challenging
 41 problem. Classical approaches [8, 9] tend to hand-design search heuristics, which requires domain
 42 knowledge and a lot of trial and error. To alleviate this problem, there has been significant development
 43 in general-purpose search heuristics based on trading-off greedy expansions and novelty-based
 44 exploration [10–13] or search problem simplifications [14–16]. These approaches alleviate some of
 45 the common pitfalls of goal-directed heuristics, but we demonstrate that if possible, it is useful to
 46 learn domain-specific heuristics that can better exploit problem structure.

47 On the other hand, learning-based methods face a set of different challenges. Firstly, the data
 48 distribution is not i.i.d., as newly encountered graph nodes depend on past heuristic values, which
 49 means that supervised learning-based methods are not directly applicable. Secondly, heuristics should
 50 run fast, with ideally constant time complexity. Otherwise, the overall asymptotic time complexity of
 51 the search procedure could be increased. Finally, as the environment (search graph) sizes increase,
 52 reinforcement learning-based heuristic learning approaches tend to perform poorly [1]. State-of-the-
 53 art imitation learning-based methods can learn useful search heuristics [1]; however, these methods
 54 still rely on feature-engineering for a specific domain and do not generally guarantee a constant time
 55 complexity with respect to graph sizes.

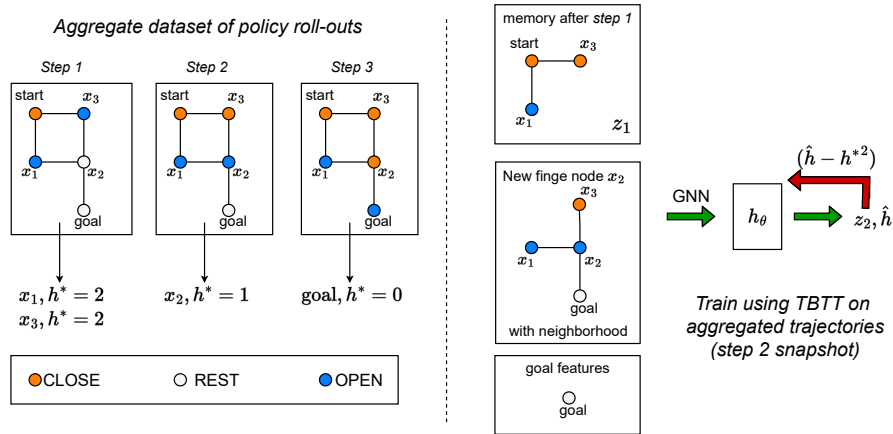


Figure 2: Main components of PHIL: On the left, using a greedy mixture policy induced by the current version of our parameterized heuristic h_θ and an oracle heuristic h^* (i.e., a heuristic that correctly determines distances between nodes), we roll-out a search trajectory from the start node to the goal node. Each trajectory step contains a set of newly added fringe nodes with bounded random subsets of their 1-hop neighborhoods and their oracle (h^*) distances to the goal node. Trajectories are aggregated throughout the training procedure. On the right, we use truncated backpropagation through time on each collected trajectory to train h_θ , where \hat{h} is the predicted distance between x_2 and x_g , and z_2 is the updated state of the memory. Here, the memory captures the embedding of the graph visited so far.

56 In this paper, we propose *Path Heuristic with Imitation Learning* (PHIL), a framework that extends
 57 the recent imitation learning-based heuristic search paradigm with a learnable *explored graph memory*.
 58 This means that PHIL learns a representation that allows it to capture the structure of the so far
 59 explored graph, so that it can then better select what node to explore next (Figure 2). We train our
 60 approach to predict the node-to-goal distances (h^* in Figure 2) of graph nodes during search. To train

61 our memory module, which captures the explored graph, we use truncated backpropagation through
 62 time (TBTT) [17], where we utilize ground-truth node-to-goal distances as a supervision signal at
 63 each search step. Our TBTT procedure is embedded within an adaptation of the AggreVaTe imitation
 64 learning algorithm [18]. PHIL also includes a *specialized graph neural network architecture*, which
 65 allows us to apply PHIL to diverse graphs from different Fdomains.

66 We evaluate PHIL on standard benchmark heuristic learning datasets (Section 5.1), diverse graph-
 67 based datasets from different domains (Section 5.2), and practical UAV flight use cases (Section 5.3).
 68 Experiments demonstrate that PHIL outperforms state-of-the-art heuristic learning methods up to
 69 $4\times$. Further, PHIL performs within 4.9% of an oracle in indoor drone planning scenarios, which is
 70 up to a 21.5% reduction compared with commonly used approaches. In practice, our contributions
 71 enable practitioners to quickly extract useful search heuristics from their graph datasets without any
 72 hand-engineering.

73 2 Preliminaries

74 **Graph search.** Suppose that we are given an unweighted connected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V}
 75 is a set of nodes, and \mathcal{E} a corresponding set of edges. Further suppose that each node $i \in \mathcal{V}$ has
 76 corresponding features $x_i \in \mathbb{R}^{D_v}$, and each edge $(i, j) \in \mathcal{E}$ has features $e_{ij} \in \mathbb{R}^{D_e}$. Assume that we
 77 are also given a start node $v_s \in \mathcal{V}$ and a goal node $v_g \in \mathcal{V}$. At any stage of our search algorithm,
 78 we can partition the nodes of our graph into three sets as $\mathcal{V} = \text{CLOSE} \cup \text{OPEN} \cup \text{REST}$, where
 79 CLOSE are the nodes already explored, OPEN are candidate nodes for exploration (i.e., all nodes
 80 connected to any node in CLOSE, but not yet in CLOSE), and REST is the rest of the graph. Each
 81 *expansion* moves a node from OPEN to CLOSE, and adds the neighbors of the given node from
 82 REST to OPEN. We call the set of newly added fringe nodes \mathcal{V}_{new} at each search step. At the start
 83 of the search procedure, $\text{CLOSE} = \{v_s\}$ and we expand the nodes until v_g is encountered (i.e., until
 84 $v_g \in \text{CLOSE}$).

85 **Greedy best-first search.** We can perform *greedy best-first search* using a greedy fringe expansion
 86 policy, such that we always expand the node $v \in \text{OPEN}$ that minimizes $h(v, v_g)$. Here, $h : \mathcal{V} \times \mathcal{V} \rightarrow$
 87 \mathbb{R} is a tailored heuristic function for a given use case. In our work, we are interested in learning a
 88 function h that predicts shortest path lengths, this way minimizing $|\text{CLOSE}|$ in a *greedy best-first*
 89 *search* regime.

90 **Imitation of perfect heuristics.** Partially observable Markov decision processes (POMDPs) are
 91 a suitable framework to describe the problem of learning search heuristics [1]. We can have $s =$
 92 (CLOSE, OPEN, REST) as our state, an action $a \in \mathcal{A}$ corresponds to moving a node from OPEN to
 93 CLOSE, and the observations $o \in \mathcal{O}$ are the features of newly included nodes in OPEN. Note that
 94 one could consider an MDP framework to learn heuristics, but the time complexity of operating on the
 95 whole state is in most cases prohibitive. We also define a history $\psi \in \Psi$ as a sequence of observations
 96 $\psi = o_1, o_2, o_3, \dots$. Our work leverages the observation that using a heuristic function during greedy
 97 best-first search that correctly determines the length of the shortest path between fringe nodes and the
 98 goal node will also yield minimal $|\text{CLOSE}|$. For training, we adopt a perfect heuristic h^* , similar to
 99 [1], which has full information about s during search. Such oracle can provide ground-truth distances
 100 $h^*(s, v, v_g)$, where $v \in \text{OPEN}$. To conclude, we define a *greedy best-first search policy* π_θ that uses
 101 a parameterized heuristic h_θ to expand nodes from OPEN with minimal heuristic values. One could
 102 also directly use a POMDP solver for the above-described problem, but this approach is usually
 103 infeasible due to the dimensionality of the search state [19].

104 3 Related Work

105 **General purpose heuristic design.** There has been significant research in designing general-purpose
 106 heuristics for speeding up satisficing planning. The first set of approaches are based on simplifying
 107 the search problem for example using landmark heuristics [14, 16]. The next set of approaches aim
 108 to include novelty-based exploration in greedy best-first search [10–13]. The latter set of approaches
 109 showed state-of-the-art performance (best-first width search [12, 13], BFWS) in numerous settings.
 110 We show that in domains where data is available, it can be more effective to incorporate a learned
 111 heuristic into a greedy best-first search procedure.

112 **Learning heuristic search.** There have been numerous previous works that attempt to learn search
 113 heuristics: Arfae *et al.* [20] propose to improve heuristics iteratively, Virseda *et al.* [21] learn to
 114 combine heuristics to estimate graph node distances, Wilt *et al.* [22] and Garrett *et al.* [23] propose
 115 to learn node rankings, Thayer *et al.* [24] suggest to infer heuristics during a search, and Kim *et al.*
 116 [25] train a neural network to predict graph node distances. These methods generally do not consider
 117 the non-i.i.d. nature of heuristic search. Further, Bhardwaj *et al.* [1] propose SAIL, where heuristic
 118 learning is framed as an imitation learning problem with cost-to-go oracles. The SAIL heuristic uses
 119 hand-designed features tailored for obstacle avoidance, with a linear time-complexity in the number
 120 of explored grid nodes found to be colliding with an obstacle. Feature-engineering becomes more
 121 difficult as we attempt to learn heuristics on diverse graphs such as ones seen in Section 5.2, where
 122 we may need expert knowledge. Further, heuristics that do not have a constant time complexity in the
 123 size of the graph [1, 26–29] generally scale poorly with graph size and hence have constrained use
 124 cases. Recent approaches to learning heuristics include Retro* [2] by Chen *et al.*, where a heuristic
 125 is learned in the context of AND-OR search trees for chemical retrosynthetic planning. Our work
 126 focuses on a more general graph setting.

127 There has been significant progress on learning heuristics for NP-hard combinatorial optimization
 128 problems [30–32]. Still, these heuristic learning methods, due to their time complexities, are
 129 impractical for the application in polynomial-time search problems, on which this work focuses.

130 **Learning general purpose search.** Learning general search policies is a very well-studied research
 131 area with a rich set of developments and applications. These include Monte Carlo Tree Search
 132 methods [33, 34], implicit planning methods [35–37], and imagination-based planning approaches
 133 [38, 39]. Learning search heuristics can be seen as a special case of *general purpose search*, where
 134 the search problem is treated as a partially observable Markov decision process with restricted action
 135 evaluation (see Section 4), and with models running in $\mathcal{O}(1)$ to remain competitive time-complexity-
 136 wise on problems where best-first search performs well. *General purpose search* methods do not
 137 take into account the above-mentioned constraints, which motivates the development of tailored
 138 approaches for learning heuristics [1, 2].

139 **Imitation learning.** Our approach builds on prior work in imitation learning (IL) with cost-to-go
 140 oracles. Cost-to-go oracles have been incorporated in the context of IL in methods such as SEARN
 141 [40], AggreVaTe [18], LOLS [41], AggreVaTeD [42], DART [43], and THOR [44]. SAIL [1] presents
 142 an AggreVaTe-based algorithm for learning heuristic search. We extend SAIL by incorporating a
 143 recurrent Q -like function, in which sense our algorithm more closely resembles AggreVaTeD by Sun
 144 *et al.* [42]. While a recurrent policy can be easily incorporated in AggreVaTeD, we cannot use a
 145 policy to evaluate actions. This is due to the fact that we would either have to evaluate all actions in a
 146 state, which is computationally infeasible, or we would have to give up on taking actions that are not
 147 in the most recent version of the search fringe, which would degrade the performance (see Section 4).

148 4 Path Heuristic with Imitation Learning

149 **Training objective.** With the aim of minimizing $|\text{CLOSE}|$ after search, our goal is to train a
 150 parameterized heuristic function $h_\theta : \Psi \times \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ to predict ground-truth node distances h^*
 151 and use h_θ within a greedy best-first policy π_θ at test time. More specifically, we assume access
 152 to a distribution over graphs P_G , a start-goal node distribution $P_{v_s g}(\cdot | \mathcal{G})$, and a time horizon T .
 153 Moreover, we assume a joint state-history distribution $s, \psi \sim P_s(\cdot | \mathcal{G}, t, \pi_\theta, v_s, v_g)$, where P_s
 154 represents the probability our search being in state s , at time $0 \leq t \leq T$ on graph \mathcal{G} with pathfinding
 155 problem (v_s, v_g) , with a greedy best-first search policy π_θ using heuristic h_θ . Hence, our goal can be
 156 summarized as minimizing the following objective:

$$\mathcal{L}(\theta) = \mathbb{E}_{\substack{\mathcal{G} \sim P_G, \\ (v_s, v_g) \sim P_{v_s g}, \\ t \sim \mathcal{U}(0, \dots, T), \\ s, \psi \sim P_s}} \left[\frac{1}{|\text{OPEN}|} \sum_{v \in \text{OPEN}} (h^*(s, v, v_g) - h_\theta(\psi, v, v_g))^2 \right] \quad (1)$$

157 Before we describe the algorithm that can be used to minimize \mathcal{L} , we rewrite h_θ to include a memory
 158 digest component (z_t) , which represents an embedding of ψ at time step t . Hence, h_θ becomes
 159 $h_\theta : \mathbb{R}^d \times \mathcal{O} \times \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$, where d is the dimensionality of our memory’s embedding space. As
 160 opposed to previous methods [1], z_t allows us to *automatically extract* relevant features for heuristic

Algorithm 1: PHIL— Sequential Heuristic Training

```

Obtain hyperparameters  $T, \beta_0, N, m, t_\tau$ ;
Initialize  $\mathcal{D} \leftarrow \emptyset, h_{\theta_1}$ ;
for  $i = 1, \dots, N$  do
  Sample  $\mathcal{G} \sim P_{\mathcal{G}}$ ;
  Sample  $v_s, v_g \sim P_{v_{sg}}$ ;
  Set  $\beta \leftarrow \beta_0^i$ ;
  Set mixture policy  $\pi_{mix} \leftarrow (1 - \beta) * \pi_{\theta_i} + \beta * \pi^*$ ;
  Collect  $m$  trajectories  $\tau_{ij}$  as follows;
  for  $j = 1, \dots, m$  do
    Sample  $t \sim \mathcal{U}(0, \dots, T - t_\tau)$ ;
    Roll-in  $t$  time steps of  $\pi_{\theta_i}$  to obtain  $z_t$  and new state  $s_t = (\text{CLOSE}^0, \text{OPEN}^0, \text{REST}^0)$ ;
    Roll-out trajectory  $\tau_{ij}$  as follows;
    for  $k = 1, \dots, t_\tau$  do
      Update  $s_{t+k-1}$  using  $\pi_{mix}$  to get new state  $s_{t+k}$  and new fringe state  $\text{OPEN}^k$ ;
      Obtain new fringe nodes  $\mathcal{V}_{new} = \text{OPEN}^k \setminus \text{OPEN}^{k-1}$ ;
      Update trajectory  $\tau_{ij} \leftarrow \tau_{ij} \cup \{(\mathcal{V}_{new}, h^*(s_{t+k}, \mathcal{V}_{new}, v_g))\}$ ;
    Update dataset  $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\tau_{ij}, z_t)\}$  or  $\mathcal{D} \cup \{(\tau_{ij}, 0)\}$ ;
  Train  $h_{\theta_i}$  using TBTT on each  $\tau \in \mathcal{D}$  to get  $h_{\theta_{i+1}}$ ;
return best performing  $h_{\theta_i}$  on validation;

```

161 computations and concurrently *reduce the computational complexity* of the heuristic function. Further,
 162 as shown in [1], if we would use h_θ to evaluate all actions in a state (i.e., recalculate the heuristic
 163 values of all nodes in OPEN), we would need a squared reduction in the number of expanded nodes
 164 compared with BFS for PHIL to bring performance benefits over BFS, which however may not
 165 be possible on all datasets. Hence, we constrain the heuristic only to evaluate new OPEN nodes
 166 which we obtain after moving a node to CLOSE, calling the set of new fringe nodes \mathcal{V}_{new} after each
 167 expansion. In practice, the policy π_θ yields an algorithm equivalent to greedy best-first search, with
 168 the heuristic function replaced by h_θ .

169 4.1 Learning algorithm & architecture

170 **Imitation learning algorithm.** In Algorithm 1, we present the pseudo-code of the IL algorithm used
 171 to train our heuristic models (Figure 3). The high-level idea of our algorithm is that we aggregate
 172 trajectories of search traces (i.e., sequences of new fringe nodes) and use truncated backpropagation
 173 through time to optimize h_θ after each data-collection step. In particular, after sampling a graph
 174 \mathcal{G} and a search problem v_s, v_g , we use our greedy learned policy π_θ induced by h_θ to roll-in for
 175 $t \sim \mathcal{U}(0, \dots, T - t_\tau)$ expansions, where T is the episode time horizon, and t_τ is the roll-out length.
 176 From our roll-in, we obtain a new state $s = (\text{CLOSE}^0, \text{OPEN}^0, \text{REST}^0)$, and an initial memory state
 177 z_t . After our roll-in, we roll-out for t_τ steps using our mixture policy π_{mix} , which is obtained by
 178 probabilistically blending π_θ and the greedy best-first policy induced by the oracle heuristic π^* . In
 179 a roll-out, we collect sequences of new fringe nodes, together with their ground-truth distances to
 180 the goal v_g , given by h^* . Once the roll-out is complete, we append the obtained trajectory and the
 181 initial state for the following optimization using backpropagation through time. Further analysis on
 182 the trade-offs between using rolled-in states z_t or zeroed-out states for training can be found in the
 183 supplementary material.

184 Note that we could also use supervised learning-based approaches to sample a fixed dataset of $(v_s,$
 185 $v_g, h^*(s, v_s, v_g))$ 3-tuples and train a model to predict node distances conditioned on their features.
 186 However, our experiments in Section 5 demonstrate that ignoring the non-i.i.d. nature of heuristic
 187 search negatively impacts model performance, with supervised learning-based methods performing
 188 up to $40\times$ worse.

189 **Recurrent GNN architecture.** In each forward pass, h_θ obtains a set of new fringe nodes \mathcal{V}_{new} , the
 190 goal node v_g , and the memory z_t at time step t . We represent each node in \mathcal{V}_{new} using its features
 191 $x_i \in \mathbb{R}^{D_v}$, and likewise the goal node v_g using its features $x_g \in \mathbb{R}^{D_v}$. Further, for each $i \in \mathcal{V}_{new}$, we
 192 uniformly sample an $n \in \mathbb{N}_{\geq 0}$ bounded set of nodes present in the 1-hop neighborhood of i , calling

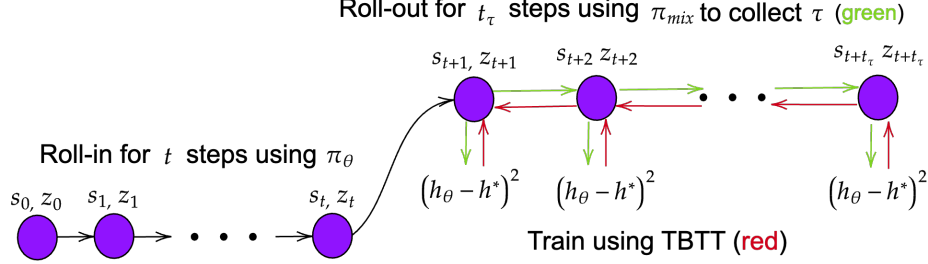


Figure 3: This figure demonstrates the core idea behind our IL algorithm. We present the roll-in phase on the left-hand side, where our policy is rolled in for t steps to obtain state s_t and embedding z_t . On the right-hand side, we show the trajectory collection and training steps, where we aggregate the trajectory for downstream training (green) and use truncated backpropagation through time on the collected dataset (red).

193 this set \mathcal{N}_i , with $|\mathcal{N}_i| \leq n$. This sampling step produces a set of neighboring node features, where
 194 each $j \in \mathcal{N}_i$ has features $x_j \in \mathbb{R}^{D_v}$, and corresponding edge features $e_{ij} \in \mathbb{R}^{D_e}$.

195 **h_θ forward pass.** Algorithm 2 presents
 196 a single forward pass of h_θ . The forward
 197 pass outputs predicted distances of the new
 198 fringe nodes to the goal \hat{h}_i , together with an
 199 updated memory digest z_{t+1} . In Algorithm
 200 2, $f, \phi, \gamma, \text{GRU}$ [45], MLP are each param-
 201 eterised differentiable functions, with ϕ, γ
 202 representing the *update* and *message* func-
 203 tions [46] of a graph neural network, re-
 204 spectively.

205 In our forward pass, using the function f ,
 206 we first project x_i, x_j into a node embed-
 207 ding space, together with the goal features x_g , and their Euclidean (D_{EUC}) and cosine distances
 208 (D_{COS}). After that, using a 1-layer GNN, we perform a single convolution over each x_i and the
 209 corresponding neighborhood \mathcal{N}_i , to obtain g_i . The specific GNN choice is a design decision left
 210 to the practitioner, and further analysis of GNN choices can be found in Appendix D. Our graph
 211 convolution processing step allows us to easily incorporate edge features and work with variable sizes
 212 of \mathcal{N}_i . After the graph convolution, we apply the GRU module over each embedding g_i to obtain
 213 hidden states $z_{i,t+1}$, and new embeddings g'_i . We compute the sample mean of $z_{i,t+1}$ for each node
 214 $i \in \mathcal{V}_{new}$ to obtain a new hidden state z_{t+1} , and process g'_i with x_g using an MLP to compute the
 215 distances between the graph nodes.

216 **Permutation invariant \mathcal{V}_{new} embedding.** There is a trade-off between processing new fringe nodes
 217 in batch, as in Algorithm 2, and processing them sequentially. Namely, when we process the nodes in
 218 batch, we do not use the in-batch observations to predict batch node values, which means that z_t
 219 is slightly outdated. On the other hand, in PHIL, batch processing allows us to compute the heuristic
 220 values of all $v \in \mathcal{V}_{new}$ in parallel on a GPU and preserves the memory’s permutation invariance
 221 with respect to nodes in \mathcal{V}_{new} . That is, because our observations are nodes & edges of a graph,
 222 the respective observation ordering usually does not contain inductive biases useful for predictions,
 223 which means that we can apply a permutation invariant operator such as the mean of all new states
 224 $z_{i,t+1}$ to obtain an aggregated updated state. This approach provides additional scalability as we can
 225 process values in parallel and PHIL does not have to infer permutation invariance in \mathcal{V}_{new} from data.

226 **Runtime complexity.** Since $\forall i \in \mathcal{V}_{new} : |\mathcal{N}_i| \leq n$, Algorithm 2 together with neighborhood
 227 sampling runs in up to $nc_1 + (n+1)c_2$ operations per each node $i \in \mathcal{V}_{new}$, which is $\mathcal{O}(1)$ with
 228 respect to the size of the graph. Here, c_1 is the maximal number of operations associated with
 229 evaluating a node, such as performing robot collision checks in dynamically constructed graphs, and
 230 c_2 is the maximal count of total model operations (e.g., f & γ operations) on the node set $\{i\} \cup \mathcal{N}_i$. In
 231 general, we expect to learn a better search heuristic with increasing n (see Appendix D for ablations),

Algorithm 2: Heuristic func. (h_θ) forward pass

Obtain $x_i, x_j, e_{ij}, x_g, z_t$;
 $x_i \leftarrow f(x_i, x_g, D_{EUC}(x_i, x_g), D_{COS}(x_i, x_g));$
 $x_j \leftarrow f(x_j, x_g, D_{EUC}(x_j, x_g), D_{COS}(x_j, x_g));$
 $g_i \leftarrow \phi(x_i, \oplus_{j \in \mathcal{N}_i} \gamma(x_i, x_j, e_{ij}));$
 $g'_i, z_{i,t+1} \leftarrow \text{GRU}(g_i, z_t);$
 $z_{t+1} \leftarrow \bar{z}_{i,t+1};$
 $\hat{h}_i \leftarrow \text{MLP}(g'_i, x_g);$
return $\hat{h}_i, z_{t+1};$

232 but in some use cases, c_1 may dominate overall complexity, which means the hyperparameter n is
 233 helpful for practitioners to tune trade-offs between constant factors and search effort minimization.

234 5 Experiments

235 In our experiments, we evaluate PHIL both on benchmark heuristic learning datasets [1] (Section
 236 5.1) as well on a diverse set of graph datasets (Section 5.2). Finally, we show that PHIL can be
 237 applied to efficient planning in the context of drone flight (Section 5.3). Our main goal is to assess
 238 how PHIL compares to baseline methods in terms of necessary expansions before the goal node is
 239 reached. Please refer to the supplementary material for information about baselines, an ablation study,
 240 and additional experiment details.

241 5.1 Heuristic search in grids

242 In Section 5.1, we evaluate PHIL on 8, 200×200 8-connected grid
 243 graph-based datasets by Bhardwaj *et al.* [1]. These datasets present
 244 challenging obstacle configurations for naive greedy planning heuristics,
 245 especially when v_s is in the bottom-left of the grid, and v_g in the
 246 top-right. Each dataset contains 200 training graphs, 70 validation
 247 graphs, and 100 test graphs. Example graphs from each dataset can
 248 be found in Table 1.

249 We train PHIL with a hyperparameter configuration of $T = 128$,
 250 $t_\tau = 32$, $\beta_0 = 0.7$, $n = 8$, and using rolled-in z_t states as initial
 251 states for training. We use a 3-layer MLP of width 128 with
 252 *LeakyReLU* activations, followed by a DeeperGCN [47] graph con-
 253 volution with *softmax* aggregation. Our memory’s embedding dimen-
 254 sionality is 64. See Appendix C for an overview of our baselines
 255 and datasets.

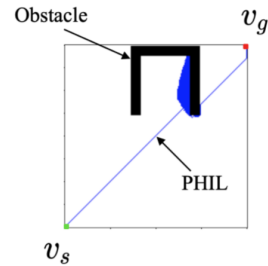


Figure 4: Example of PHIL escaping local search minima.

Dataset	Graph Examples	SAIL	SL	CEM	QL	h_{euc}	h_{man}	A*	MHA*	BFWS	Neural A*	PHIL
Alternating gaps		0.039	0.432	0.042	1.000	1.000	1.000	1.000	1.000	0.34	0.546	0.024
Single Bugtrap		0.158	0.214	0.057	1.000	0.184	0.192	1.000	0.286	0.099	0.394	0.077
Shifting gaps		0.104	0.464	1.000	1.000	0.506	0.589	1.000	0.804	0.206	0.563	0.027
Forest		0.036	0.043	0.048	0.121	0.041	0.043	1.000	0.075	0.039	0.399	0.027
Bugtrap+Forest		0.147	0.384	0.182	1.000	0.410	0.337	1.000	3.177	0.149	0.651	0.135
Gaps+Forest		0.221	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.401	0.580	0.039
Mazes		0.103	0.238	0.479	0.399	0.185	0.171	1.000	0.279	0.095	1.000	0.069
Multiple Bugtraps		0.479	0.480	1.000	0.835	0.648	0.617	1.000	0.876	0.169	0.331	0.136

Table 1: The number of expanded graph nodes of PHIL with respect to SAIL. We can observe that out of all baselines, SAIL performs best. PHIL outperforms SAIL by 58.5% on average over all datasets, with a maximal search effort reduction of 82.3% in the *Gaps+Forest* dataset.

256 **Discussion.** As we can see in Table 1, PHIL outperforms the best baseline (SAIL) on all datasets,
 257 with an average reduction of explored nodes before v_g is found of 58.5%. Qualitatively, observing
 258 Figure 5, we can attribute these results to PHIL’s ability to **reduce the redundancy in explored**
 259 **nodes** during a search. Further, PHIL is also **capable of escaping local minima**, which is illustrated
 260 in Figure 4. However, note that we occasionally observe failure cases in practice, where PHIL gets
 261 stuck in a bug trap-like structure. We discuss possible remedies and opportunities for future work in
 262 the supplementary material.

263 **Runtime & convergence speed.** PHIL converges in up to $N = 36$ iterations, with $m = 1$, $t_\tau = 32$
 264 (i.e., after observing less than $N * t_\tau * \max(|\mathcal{V}_{new}|) \approx 9,216$ shortest path distances, where we
 265 take $\max(|\mathcal{V}_{new}|) = 8$ as the maximal size of \mathcal{V}_{new}). According to figures reported in [1], this is
 266 approximately $5 \times$ less data than it takes for SAIL to converge.

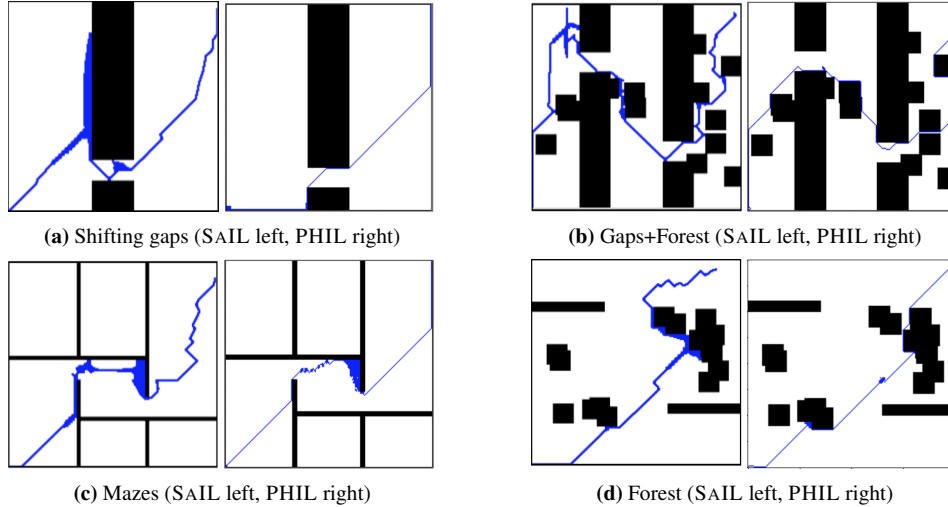


Figure 5: In each image pair of this figure, we provide a qualitative comparison with the SAIL method. In particular, we show comparisons on the *Shifting gaps*, *Gaps+Forest*, *Mazes*, and *Forest* datasets. We can observe that PHIL (right) learns the appropriate heuristics for the given dataset and makes fewer redundant expansions than SAIL (left).

267 5.2 Search in real-life graphs of different structures

268 In this experiment, our goal is to demonstrate the general applicability of PHIL to various graphs. We
 269 train PHIL on 4 different groups of graph datasets: citation networks, biological networks, abstract
 270 syntax trees (ASTs), and road networks. We use the same graph for citation networks and road
 271 networks for training and evaluation, and we use 100 random v_s, v_g pairs for testing. In the case of
 272 biological networks and ASTs, we usually have train/validation/test splits of 80/10/10, and in the
 273 case of the OGB [48] datasets, we use the provided splits.

	Dataset	$ \mathcal{D} $	$ \mathcal{V} $	$ \mathcal{E} $	SL	A*	h_{euc}	BFS	SAIL	BFWS	PHIL
Citation Networks	Cora (Sen <i>et al.</i> [49])	1	2,708	5,429	2.201	2.067	1.000	4.001	0.669	1.378	0.475
	PubMed (Sen <i>et al.</i> [49])	1	19,717	44,338	2.157	2.983	1.000	3.853	1.196	1.000	0.745
	CiteSeer (Sen <i>et al.</i> [49])	1	3,327	4,732	1.636	1.487	1.000	2.190	1.062	0.951	0.599
	Coauthor (cs) (Schur <i>et al.</i> [50])	1	18,333	81,894	1.571	1.069	1.000	2.820	1.941	1.026	0.835
	Coauthor (physics) (Schur <i>et al.</i> [50])	1	34,493	247,962	4.076	1.081	1.000	4.523	–	1.012	0.964
Biological Networks	OGBG-Molhiv (Hu <i>et al.</i> [48])	41,127	25.5	27.5	1.086	1.065	1.000	1.267	1.104	1.146	1.016
	PPI (Zitnik <i>et al.</i> [51])	24	2,372.67	34,113.16	0.772	0.831	1.000	5.618	1.746	3.941	0.658
	Proteins (Full) (Morris <i>et al.</i> [52])	1,113	39.06	72.82	0.995	0.997	1.000	2.645	0.891	0.966	0.831
	Enzymes (Morris <i>et al.</i> [52])	600	32.63	62.14	1.073	1.007	1.000	1.358	1.036	0.992	0.757
ASTs	OGBG-Code2 (Hu <i>et al.</i> [48])	452,741	125.2	124.2	1.196	1.013	1.000	1.267	1.029	0.817	1.219
Road Networks	OSMnx - Modena (Boeing [53])	1	29,324	38,309	2.904	3.085	1.000	3.493	1.182	0.997	0.489
	OSMnx - New York (Boeing [53])	1	54,128	89,618	39.424	36.529	1.000	63.352	1.583	1.013	0.962

Table 2: Comparison of PHIL with baseline approaches on 4 groups of datasets: citation networks, biological networks, abstract syntax trees, and road networks. "–" denotes being out of a 4-day's training time limit. We can observe that, on average across all datasets, PHIL outperforms the best baseline per dataset by 13.4%. Discounting the OGBG datasets, this number becomes 19.5%.

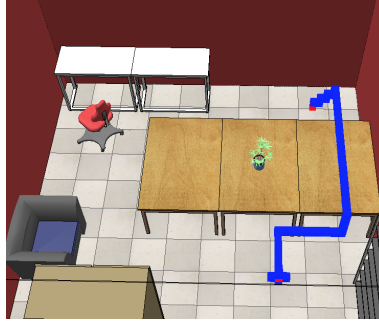
274 Similarly as in Section 5.1, our MLP has four layers of width 128 with *LeakyReLU* activations and
 275 we use a DeeperGCN [47] graph convolution with *softmax* aggregation. The utilized node and edge
 276 features are the provided features in each dataset, except for a few minor modifications which are
 277 discussed in Appendix A & Appendix C. We train an MLP of depth 5 and width 256 using supervised
 278 learning (SL) for our learning-based baseline method.

279 **Discussion.** The results presented in Table 2 suggest that PHIL can learn superior search heuristics
 280 compared with baseline methods, outperforming top baselines per dataset in terms of visited nodes
 281 during a search by 13.4% on average. Two datasets where PHIL fell short compared to other
 282 baselines are the *OGBG-Molhiv* and *OGBG-Code2* datasets. The *OGBG-Code2* dataset adopts a
 283 *project split* [54] and *OGBG-Molhiv* adopts a *scaffold split* [55], both of which ensure that graphs of
 284 different structure are present in the training & test sets. Although PHIL improved upon uninformed
 285 search (BFS) in the OGB datasets, structural graph consistency is explicitly discouraged in the

286 above-mentioned OGBG splits. Without the OGBG datasets, PHIL improves on the top baselines per
 287 dataset by 19.5% on average, and upon the Euclidean node feature heuristic (h_{euc}) by 20.4%. Note
 288 that we trained PHIL up to $N = 60$ iterations, which means that it only encountered a small subset
 289 of the pathfinding problems in the single graph setting, which means that PHIL had to generalize to
 290 learn useful heuristics. Even in Cora, the $|\mathcal{D}| = 1$ dataset with least number of nodes, PHIL observed
 291 roughly 6,000 node distances during training, which is less than 0.2% of total distances in the Cora
 292 graph.

293 5.3 Planning for drone flight

294 In our final experiment, we use PHIL to plan collision-
 295 free paths in a practical drone flight use case within an
 296 indoor environment. We built our environment using the
 297 CoppeliaSim simulator [56], and the Ivy framework [57].
 298 Figure 6 presents the environment which we refer to as
 299 *room adversarial* in Table 3. For more detail about each
 300 environment, please refer to the supplementary material.
 301 We discretize the environments into 3D grid graphs of size
 302 $50 \times 50 \times 25$, and randomly remove 5 sub-graphs of size $5 \times$
 303 5×5 both during training and testing, this way simulating
 304 real-life planning scenarios with random obstacles. The
 305 hyperparameter configuration and the specific architecture
 306 we utilize are equivalent to Section 5.1, but with $n = 4$.
 307 Likewise, the node features are 3D grid coordinates, and
 308 the baselines include supervised learning (SL), h_{euc} , A*,
 309 and BFS, similarly as in Sections 5.1, 5.2. In Table 3 we
 310 report the ratio of expanded nodes with respect to h_{euc} .



311 **Figure 6:** This figure illustrates the *room adversarial* environment with an exam-
 312 ple planning problem (red) and the ex-
 313 panded graph by PHIL (blue).

314 **Video demo.** We provide a video demonstration of PHIL running in *room adversarial*: <https://cutt.ly/eniu5ax>.

Dataset	SL	A*	h_{euc}	BFS	SAIL	BFWS	PHIL	Shortest path
Room simple	1.124	76.052	1.000	291.888	0.973	1.286	0.785	0.782
Room adversarial	2.022	67.215	1.000	238.768	0.944	1.583	0.895	0.853

315 **Table 3:** Results of PHIL in the context of planning for indoor UAV flight. PHIL outperforms
 316 all baselines in both the *room simple* and *room adversarial* environments while remaining close
 317 performance-wise to the optimal number of expansions.

318 **Discussion.** As we can observe in Table 3, PHIL outperforms all baselines in both environments.
 319 Interestingly, PHIL expands only approximately 0.3% more nodes in the simple room than least
 320 possible and 4.9% more in the adversarial room case. The same figures for the greedy method (h_{euc})
 321 are 27.8% and 17.2%, respectively. These results indicate that PHIL is capable of learning planning
 322 strategies that are close to optimal in both *simple* and *adversarial* graphs, while the performance of
 323 naive heuristics degrades.

319 5.4 Runtime Analysis

320 We summarize test run-times of different approaches in Appendix G. PHIL runs 57.9% faster than
 321 BFWS and 32.2% faster than SAIL, and not much slower than traditional A* (34.7%) and h_{man}
 322 (18.3%). Although Neural A* is 71.0% faster than PHIL due to the fact that it casts the whole search
 323 process into matrix operations on images, it cannot be employed in a generic search setting.

324 6 Conclusion

325 In our work, we consider the problem of learning to search for feasible paths in graphs
 326 efficiently. We propose a model and a training procedure to learn search heuristics that
 327 can be easily deployed across diverse graphs, with tuneable trade-off parameters between
 328 constant factors and performance. Our results demonstrate that PHIL outperforms cur-
 329 rent state-of-the-art approaches and can be applied to various graphs with practical use
 330 cases.

References

- 331
- 332 [1] Mohak Bhardwaj, Sanjiban Choudhury, and Sebastian Scherer. Learning heuristic search via
333 imitation. In *Conference on Robot Learning*, 2017. 1, 2, 3, 4, 5, 7, 13, 14, 15, 16, 19, 22
- 334 [2] Binghong Chen, Chengtao Li, Hanjun Dai, and Le Song. Retro*: Learning retrosynthetic
335 planning with neural guided a* search. In *ICML*, 2020. 1, 4
- 336 [3] Martin Gebser, Benjamin Kaufmann, Javier Romero, Ramón Otero, Torsten Schaub, and Philipp
337 Wanko. Domain-specific heuristics in answer set programming. In *AAAI*, 2013. 1
- 338 [4] Thi Thoa Mac, Cosmin Copot, Duc Trung Tran, and Robin De Keyser. Heuristic approaches in
339 robot path planning: A survey. In *Robotics and Autonomous Systems*, 2016. 1
- 340 [5] Abhishek Sharma and Keith M. Goolsbey. Identifying useful inference paths in large common-
341 sense knowledge bases by retrograde analysis. In *AAAI*, 2017. 1
- 342 [6] Cheng-Yu Yeh, Hsiang-Yuan Yeh, Carlos Roberto Arias, and Von-Wun Soo. Pathway detection
343 from protein interaction networks and gene expression data using color-coding methods and a*
344 search algorithms. In *The Scientific World booktitle*, 2012. 1
- 345 [7] Judea Pearl. *Heuristics: intelligent search strategies for computer problem solving*. 1984. 2
- 346 [8] Danish Khalidi, Dhaval Gujarathi, and Indranil Saha. T*: A heuristic search based path planning
347 algorithm for temporal logic specifications. In *ICRA*, 2020. 2
- 348 [9] Bhargav Adabala and Zlatan Ajanovic. A multi-heuristic search-based motion planning for
349 autonomous parking. In *30th International Conference on Automated Planning and Scheduling:
350 Planning and Robotics Workshop*, 2020. 2
- 351 [10] Fan Xie, Hootan Nakhost, and Martin Müller. Planning via random walk-driven local search.
352 In *Twenty-Second International Conference on Automated Planning and Scheduling*, 2012. 2, 3
- 353 [11] Fan Xie, Martin Müller, and Robert Holte. Adding local exploration to greedy best-first search
354 in satisficing planning. In *AAAI*, 2014.
- 355 [12] Nir Lipovetzky and Hector Geffner. Best-first width search: Exploration and exploitation in
356 classical planning. In *AAAI*, 2017. 3, 14
- 357 [13] Florent Teichteil-Königsbuch, Miquel Ramirez, and Nir Lipovetzky. Boundary extension
358 features for width-based planning with simulators on continuous-state domains. In *Proceedings
359 of the Twenty-Ninth International Conference on International Joint Conferences on Artificial
360 Intelligence*, 2021. 2, 3, 14
- 361 [14] Blai Bonet and Héctor Geffner. Planning as heuristic search. pages 5–33. 2, 3
- 362 [15] Lin Zhu and Robert Givan. Landmark extraction via planning graph propagation. 2003.
- 363 [16] Silvia Richter and Matthias Westphal. The lama planner: Guiding cost-based anytime planning
364 with landmarks. 2010. 2, 3
- 365 [17] Ilya Sutskever. *Training recurrent neural networks*. University of Toronto, Toronto, Canada,
366 2013. 3
- 367 [18] Stephane Ross and J Andrew Bagnell. Reinforcement and imitation learning via interactive
368 no-regret learning. In *arXiv preprint arXiv:1406.5979*, 2014. 3, 4
- 369 [19] Sanjiban Choudhury, Ashish Kapoor, Gireeja Ranade, Sebastian Scherer, and Debadepta Dey.
370 Adaptive information gathering via imitation learning. 2017. 3
- 371 [20] Shahab Jabbari Arfaee, Sandra Zilles, and Robert C Holte. Learning heuristic functions for
372 large state spaces. In *Artificial Intelligence*, 2011. 4
- 373 [21] Jes ús Virseda, Daniel Borrajo, and Vidal Alcázar. Learning heuristic functions for cost-based
374 planning. In *Planning and Learning*, 2013. 4
- 375 [22] Christopher Makoto Wilt and Wheeler Ruml. Building a heuristic for greedy search. In *SOCS*,
376 2015. 4
- 377 [23] Caelan Reed Garrett, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. Learning to rank for
378 synthesizing planning heuristics. In *IJCAI*, 2016. 4
- 379 [24] Jordan Thayer, Austin Dionne, and Wheeler Ruml. Learning inadmissible heuristics during
380 search. In *Proceedings of the International Conference on Automated Planning and Scheduling*,
381 2011. 4

- 382 [25] Soonkyum Kim and Byungchul An. Learning heuristic a*: efficient graph search using neural
383 network. In *ICRA*, 2020. 4
- 384 [26] Yuka Ariki and Takuya Narihira. Fully convolutional search heuristic learning for rapid path
385 planners. In *arXiv preprint arXiv:1908.03343*, 2019. 4
- 386 [27] Ryo Terasawa, Yuka Ariki, Takuya Narihira, Toshimitsu Tsuboi, and Kenichiro Nagasaka.
387 3d-cnn based heuristic guided task-space planner for faster motion planning. In *ICRA*, 2020.
- 388 [28] Ryo Yonetani, Tatsunori Tani, Mohammadamin Barekatin, Mai Nishimura, and Asako
389 Kanazaki. Path planning using neural a* search. In *ICML*, 2021. 14
- 390 [29] Alberto Archetti, Marco Cannici, and Matteo Matteucci. Neural weighted a*: Learning graph
391 costs and heuristics with differentiable anytime a. 2021. 4
- 392 [30] Elias B. Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilikina, and Le Song. Learning combinatorial
393 optimization algorithms over graphs. In *NeurIPS*, 2017. 4
- 394 [31] Zhuwen Li, Qifeng Chen, and Vladlen Koltun. Combinatorial optimization with graph convolu-
395 tional networks and guided tree search. In *NeurIPS*, 2018.
- 396 [32] Nikolaos Karalias and Andreas Loukas. Erdos goes neural: an unsupervised learning framework
397 for combinatorial optimization on graphs. In *NeurIPS*, 2020. 4
- 398 [33] David Silver and Joel Veness. Monte-carlo planning in large pomdps. In *NeurIPS*, 2010. 4
- 399 [34] Arthur Guez, Theophane Weber, Ioannis Antonoglou, Karen Simonyan, Oriol Vinyals, Daan
400 Wierstra, Rémi Munos, and David Silver. Learning to search with mctsnets. In *ICML*, 2018. 4
- 401 [35] Andreea Deac, Petar Veličković, Ognjen Milinković, Pierre-Luc Bacon, Jian Tang, and Mladen
402 Nikolić. Xlvin: executed latent value iteration nets. In *arXiv preprint arXiv:2010.13146*, 2020.
403 4
- 404 [36] Péter Karkus, David Hsu, and Wee Sun Lee. Qmdp-net: Deep learning for planning under
405 partial observability. In *NeurIPS*, 2017.
- 406 [37] Aviv Tamar, Sergey Levine, Pieter Abbeel, Yi Wu, and Garrett Thomas. Value iteration networks.
407 In *NeurIPS*, 2016. 4
- 408 [38] Razvan Pascanu, Yujia Li, Oriol Vinyals, Nicolas Heess, Lars Buesing, Sebastien Racanière,
409 David Reichert, Théophane Weber, Daan Wierstra, and Peter Battaglia. Learning model-based
410 planning from scratch. In *arXiv preprint arXiv:1707.06170*, 2017. 4
- 411 [39] Sébastien Racanière, Theophane Weber, David P. Reichert, Lars Buesing, Arthur Guez,
412 Danilo Jimenez Rezende, Adrià Puigdomènech Badia, Oriol Vinyals, Nicolas Heess, Yujia
413 Li, Razvan Pascanu, Peter W. Battaglia, Demis Hassabis, David Silver, and Daan Wierstra.
414 Imagination-augmented agents for deep reinforcement learning. In *NeurIPS*, 2017. 4
- 415 [40] Hal Daumé, John Langford, and Daniel Marcu. Search-based structured prediction. In *Machine
416 learning*, 2009. 4
- 417 [41] Kai-Wei Chang, Akshay Krishnamurthy, Alekh Agarwal, Hal Daumé III, and John Langford.
418 Learning to search better than your teacher. In *ICML*, 2015. 4
- 419 [42] Wen Sun, Arun Venkatraman, Geoffrey J. Gordon, Byron Boots, and J. Andrew Bagnell. Deeply
420 aggravated: differentiable imitation learning for sequential prediction. In *ICML*, 2017. 4
- 421 [43] Michael Laskey, Jonathan Lee, Roy Fox, Anca Dragan, and Ken Goldberg. Dart: Noise injection
422 for robust imitation learning. In *Conference on robot learning*, 2017. 4
- 423 [44] Wen Sun, J. Andrew Bagnell, and Byron Boots. Truncated horizon policy search: combining
424 reinforcement learning & imitation learning. In *ICLR*, 2018. 4
- 425 [45] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares,
426 Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-
427 decoder for statistical machine translation. In *EMNLP*, 2014. 6
- 428 [46] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl.
429 Neural message passing for quantum chemistry. In *ICML*, 2017. 6, 17
- 430 [47] Guohao Li, Chenxin Xiong, Ali Thabet, and Bernard Ghanem. Deepergcn: All you need to
431 train deeper gcn. In *arXiv preprint arXiv:2006.07739*, 2020. 7, 8, 17

- 432 [48] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
433 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.
434 In *NeurIPS*, 2020. 8
- 435 [49] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-
436 Rad. Collective classification in network data. In *AI magazine*, 2008. 8
- 437 [50] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann.
438 Pitfalls of graph neural network evaluation. In *arXiv preprint arXiv:1811.05868*, 2018. 8
- 439 [51] Marinka Zitnik and Jure Leskovec. Predicting multicellular function through multi-layer tissue
440 networks. In *Bioinformatics*, 2017. 8
- 441 [52] Christopher Morris, Nils M Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion
442 Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. In *arXiv
443 preprint arXiv:2007.08663*, 2020. 8
- 444 [53] Geoff Boeing. Osmnx: New methods for acquiring, constructing, analyzing, and visualizing
445 complex street networks. In *Computers, Environment and Urban Systems*, 2017. 8
- 446 [54] Miltiadis Allamanis. The adverse effects of code duplication in machine learning models
447 of code. In *ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and
448 Reflections on Programming and Software*, 2019. 8
- 449 [55] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S
450 Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine
451 learning. In *Chemical science*, 2018. 8
- 452 [56] E. Rohmer, S. P. N. Singh, and M. Freese. Coppeliassim (formerly v-rep): a versatile and scalable
453 robot simulation framework. In *IROS*, 2013. 9, 15
- 454 [57] Daniel Lenton, Fabio Pardo, Fabian Falck, Stephen James, and Ronald Clark. Ivy: Templated
455 deep learning for inter-framework portability. In *arXiv preprint arXiv:2102.02886*, 2021. 9, 15
- 456 [58] Jiaxuan You, Rex Ying, and Jure Leskovec. Position-aware graph neural networks. In *ICML*,
457 2019. 13
- 458 [59] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *ACM
459 SIGKDD*, 2016. 13
- 460 [60] Stuart Russell and Peter Norvig. *Artificial intelligence: a modern approach*. 2002. 13
- 461 [61] Sandip Aine, Siddharth Swaminathan, Venkatraman Narayanan, Victor Hwang, and Maxim
462 Likhachev. Multi-heuristic a*. In *The International booktitle of Robotics Research*, 2016. 13,
463 14
- 464 [62] Edo Cohen-Karlik, Avichai Ben David, and Amir Globerson. Regularizing towards permutation
465 invariance in recurrent models. In *NeurIPS*, 2020. 13
- 466 [63] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
467 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. In *NeurIPS
468 Deep Learning Workshop*, 2013. 14
- 469 [64] Pieter-Tjerk De Boer, Dirk P Kroese, Shie Mannor, and Reuven Y Rubinstein. A tutorial on the
470 cross-entropy method. In *Annals of operations research*, 2005. 14
- 471 [65] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
472 Bengio. Graph attention networks. In *ICLR*, 2018. 17
- 473 [66] Petar Velickovic, Rex Ying, Matilde Padovano, Raia Hadsell, and Charles Blundell. Neural
474 execution of graph algorithms. In *ICLR*, 2020. 18
- 475 [67] Petar Velickovic. Tikz. <https://github.com/PetarV-/TikZ>, last accessed on 01/6/21. 20