

---

# Deep Policy Dynamic Programming for Vehicle Routing Problems

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Routing problems are a class of combinatorial problems with many practical  
2 applications. Recently, end-to-end deep learning methods have been proposed  
3 to learn approximate solution heuristics for such problems. In contrast, classical  
4 dynamic programming (DP) algorithms guarantee optimal solutions, but scale  
5 badly with the problem size. We propose *Deep Policy Dynamic Programming*  
6 (DPDP), which aims to combine the strengths of learned neural heuristics with  
7 those of DP algorithms. DPDP prioritizes and restricts the DP state space using  
8 a policy derived from a deep neural network, which is trained to predict edges  
9 from example solutions. We evaluate our framework on the travelling salesman  
10 problem (TSP), the vehicle routing problem (VRP) and TSP with time windows  
11 (TSPTW) and show that the neural policy improves the performance of (restricted)  
12 DP algorithms, making them competitive to strong alternatives such as LKH, while  
13 also outperforming most other ‘neural approaches’ for solving TSPs, VRPs and  
14 TSPTWs with 100 nodes.

## 15 1 Introduction

16 Dynamic programming (DP) is a powerful framework for solving optimization problems by solving  
17 smaller subproblems through the principle of optimality [3]. Famous examples are Dijkstra’s  
18 algorithm [14] for the shortest route between two locations, and the classic Held-Karp algorithm for  
19 the travelling salesman problem (TSP) [23, 4]. Despite their long history, dynamic programming  
20 algorithms for vehicle routing problems (VRPs) have seen limited use in practice, primarily due to  
21 their bad scaling performance. More recently, a line of research has attempted the use of machine  
22 learning (especially deep learning) to automatically learn heuristics for solving routing problems  
23 [57, 5, 41, 29, 7]. While the results are promising, most learned heuristics are not (yet) competitive  
24 to ‘traditional’ algorithms such as LKH [24] and lack (asymptotic) guarantees on their performance.

25 In this paper, we propose *Deep Policy Dynamic Programming* (DPDP) as a framework for solving  
26 vehicle routing problems. The key of DPDP is to combine the strengths of deep learning and DP,  
27 by restricting the DP state space (the search space) using a policy derived from a neural network.  
28 In Figure 1 it can be seen how the neural network indicates promising parts of the search space  
29 (through a *heatmap* over the edges of the graph), which is then used by the DP algorithm to find a  
30 good solution. DPDP is more powerful than some related ideas [62, 52, 61, 6, 34] as it combines  
31 supervised training of a large neural network with just a *single* model evaluation at test time, to enable  
32 running a large scale guided search using DP. The DP framework is flexible as it can model a variety  
33 of realistic routing problems with difficult practical constraints [20]. We illustrate this by testing  
34 DPDP on the TSP, the capacitated VRP and the TSP with (hard) time window constraints (TSPTW).

35 In more detail, the starting point of our proposed approach is a *restricted dynamic programming*  
36 algorithm [20]. Such an algorithm heuristically reduces the search space by retaining only the  $B$  most

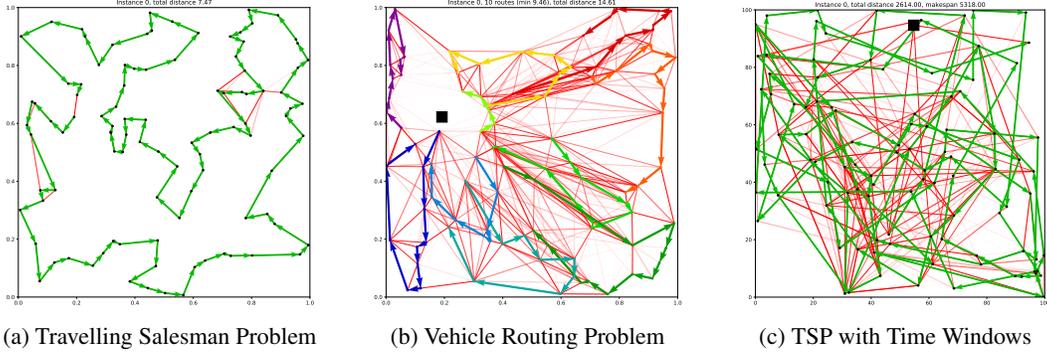


Figure 1: Heatmap predictions (red) and solutions (colored) by DPDP (VRP depot edges omitted).

37 promising solutions per iteration. The selection process is very important as it defines the part of the  
 38 DP state space considered and, thus, the quality of the solution found (see Fig. 2). Instead of manually  
 39 defining a selection criterion, DPDP defines it using a (sparse) heatmap of promising route segments  
 40 obtained by pre-processing the problem instance using a (deep) graph neural network (GNN) [26].  
 41 This approach is reminiscent of neural branching policies for branch-and-bound algorithms [19, 40].

42 In this work, we thus aim for a ‘neural boost’ of DP algorithms, by using a graph neural network  
 43 for scoring partial solutions. Prior work on ‘neural’ vehicle routing has focused on auto-regressive  
 44 models [57, 5, 13, 29], but they have high computational cost when combined with (any form of)  
 45 search, as the model needs to be evaluated for each partial solution considered. Instead, we use (for  
 46 TSP) and adapt (for VRP and TSPTW) a model to predict a heatmap indicating promising edges [26],  
 47 and define the *score* of a partial solution as the ‘heat’ of the edges it contains (plus an estimate of the  
 48 ‘heat-to-go’ or *potential* of the solution). As the neural network only needs to be evaluated *once*  
 49 for each instance, this enables a *much larger search* (defined by  $B$ ), making a good trade-off between  
 50 quality and computational cost. Additionally, we can apply a threshold to the heatmap to define a  
 51 sparse graph on which to run the DP algorithm, reducing the runtime by eliminating many solutions.

52 Figure 2 illustrates the overall DPDP algorithm. In Section 4, we show that DPDP significantly  
 53 improves over ‘classic’ restricted DP algorithms (with the same  $B$ ). Additionally, we show that  
 54 DPDP outperforms most other ‘neural’ approaches for TSP, VRP and TSPTW and is competitive  
 55 with the highly-optimized LKH solver [24] for VRP, while achieving similar results much faster for  
 56 TSP and TSPTW. For TSPTW, DPDP also outperforms the best open-source solver we could find  
 57 [10], illustrating the power of DPDP to handle difficult hard constraints (time windows).

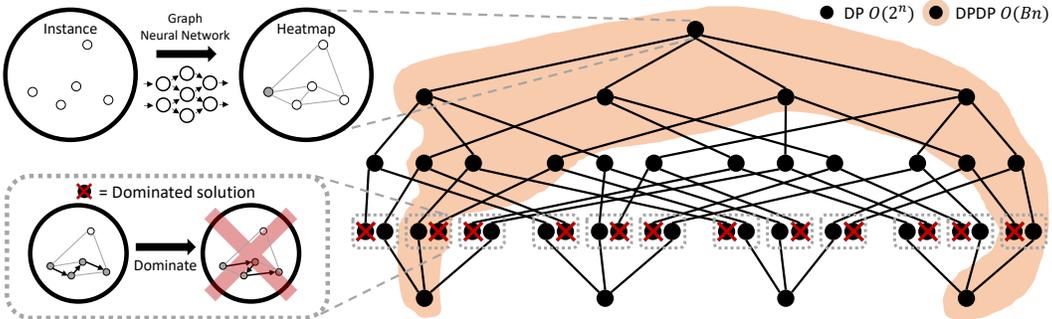


Figure 2: Deep Policy Dynamic Programming for the TSP. A GNN creates a (sparse) heatmap indicating promising edges, after which a tour is constructed using forward dynamic programming. In each step, at most  $B$  solutions are expanded according to the heatmap policy, restricting the size of the search space. Partial solutions are dominated by shorter (lower cost) solutions with the same DP state: the same nodes visited (marked grey) and current node (indicated by dashed rectangles).

## 58 2 Related work

59 DP has a long history as an exact solution method for routing problems [31, 50], e.g. for the TSP  
60 with time windows [15] and precedence constraints [39], but typically limited to small problems only,  
61 due to the curse of dimensionality. Restricted DP (with heuristic policies) has been used to address,  
62 e.g., the time dependent TSP [37], and has been generalized into a flexible framework for VRPs with  
63 different types of practical constraints [20]. DP approaches have also been shown to be useful in  
64 settings with difficult practical issues such as time-dependent travel times and driving regulations [28]  
65 or stochastic demands [42]. For a thorough investigation of modelling choices of DP for routing (and  
66 scheduling), see [53]. For sparse graphs, alternative, but less flexible, formulations can be used [8].

67 Despite the flexibility, constructive DP methods have not gained much popularity compared to  
68 heuristic search approaches such as Ruin and Recreate [47], Adaptive Large Neighborhood Search  
69 [46], LKH [24] or FILO [1]. While highly effective, these methods are limited in their flexibility as  
70 special operators need to be engineered for different types of problems. While restricted DP was  
71 shown to have superior performance on *realistic* VRPs with many constraints [20], the performance  
72 gap of around 10% for standard (benchmark) VRPs (with time windows) is too large to popularize  
73 the restricted dynamic programming approach. We argue that the missing ingredient for restricted  
74 dynamic programming is the availability of a strong but computationally cheap policy for selecting  
75 which solutions should be considered, which is the motivation behind DPDP.

76 In the machine learning community, recent advances have significantly improved deep neural networks  
77 (DNNs) to perform tasks such as image classification and machine translation [32]. After the first  
78 deep learning model was trained (using example solutions) to construct TSP tours [57], many  
79 improvements have been proposed, e.g. different training strategies such as reinforcement learning  
80 (RL) [5, 27, 12, 30] and model architectures, which enabled the same idea to be used for other  
81 routing problems [41, 29, 13, 45, 16, 60]. Most constructive neural methods are *auto-regressive*,  
82 evaluating the model many times to predict one node at the time, but other works have considered  
83 predicting a ‘heatmap’ of promising edges *at once* [43, 26, 17], which allows a tour to be constructed  
84 (using sampling or beam search) without further evaluating the model. An alternative direction is  
85 ‘learning to search’, where a neural network is used to guide a search procedure such as local search  
86 [7, 35, 18, 59, 25]. Some works have attempted scaling to larger instances beyond 100 nodes, which  
87 remains challenging [36, 17]. The combination of machine learning with DP has been proposed in  
88 limited settings [62, 52, 61]. Most related to our approach, a DP algorithm for TSPTW, guided by an  
89 RL agent, was implemented using an existing solver [6] and a neural network predicting edges has  
90 been combined with tree search [34] and local search for maximum independent set (MIS). For a  
91 wider view on machine learning for routing problems and combinatorial optimization, see [38, 54].

## 92 3 Deep Policy Dynamic Programming

93 DPDP uses an existing graph neural network [26] (which we modify for VRP and TSPTW) to predict  
94 a heatmap of promising edges, which is used to derive the policy for scoring partial solutions in  
95 the DP algorithm. The DP algorithm starts with a *beam* of a single initial (empty) solution. It then  
96 proceeds by iterating the following steps: (1) all solutions on the beam are expanded, (2) dominated  
97 solutions are removed for each *DP state*, (3) the  $B$  best solutions according to the scoring policy  
98 define the beam for the next iteration. These steps are illustrated in Fig. 2. The objective function is  
99 used to select the best solution from the final beam. The resulting algorithm is a *beam search* over the  
100 *DP state space* (which is *not* a ‘standard beam search’ over the *solution space*!) and we call  $B$  the  
101 *beam size*. DPDP is asymptotically optimal as using  $B = n \cdot 2^n$  for a TSP with  $n$  nodes guarantees  
102 optimal results, but choosing smaller  $B$  allows to trade performance for computational cost.

103 DPDP is a generic framework that can be applied to different problems, by defining the following  
104 ingredients: (1) the **state variables** to track while constructing solutions, (2) the **initial solution**,  
105 (3) **feasible actions** to expand solutions, (4) rules to define **dominated solutions** and (5) a **scoring**  
106 **policy** for selecting the  $B$  solutions to keep. A solution is always (uniquely) defined as a sequence of  
107 actions, which allows the DP algorithm to construct the final solution by backtracking. In the next  
108 sections, we define these ingredients for the TSP, VRP and TSPTW.

### 109 3.1 Travelling Salesman Problem

110 We implement DPDP for Euclidean TSPs with  $n$  nodes on a (sparse) graph, where the cost for edge  
 111  $(i, j)$  is given by  $c_{ij}$ , the Euclidean distance between the coordinates of nodes  $i$  and  $j$ .

112 For each partial solution, defined by a sequence of actions  $\mathbf{a}$ , the **state variables** are  $\text{cost}(\mathbf{a})$ , the  
 113 total *cost* (distance),  $\text{current}(\mathbf{a})$ , the current node, and  $\text{visited}(\mathbf{a})$ , the set of visited nodes (including  
 114 the start node). Without loss of generality, we let 0 be the start node, so we initialize the beam at step  
 115  $t = 0$  with the empty **initial solution** with  $\text{cost}(\mathbf{a}) = 0$ ,  $\text{current}(\mathbf{a}) = 0$  and  $\text{visited}(\mathbf{a}) = \{0\}$ . At  
 116 step  $t$ , the action  $a_t \in \{0, \dots, n - 1\}$  indicates the next node to visit, and is a **feasible action** for a  
 117 partial solution  $\mathbf{a} = (a_0, \dots, a_{t-1})$  if  $(a_{t-1}, a_t)$  is an edge in the graph and  $a_t \notin \text{visited}(\mathbf{a})$ , or, when  
 118 all are visited, if  $a_t = 0$  to return to the start node. When expanding the solution to  $\mathbf{a}' = (a_0, \dots, a_t)$ ,  
 119 we can compute the state variables incrementally as:

$$\text{cost}(\mathbf{a}') = \text{cost}(\mathbf{a}) + c_{\text{current}(\mathbf{a}), a_t}, \quad \text{current}(\mathbf{a}') = a_t, \quad \text{visited}(\mathbf{a}') = \text{visited}(\mathbf{a}) \cup \{a_t\}. \quad (1)$$

120 A (partial) solution  $\mathbf{a}$  is a **dominated solution** if there exists a (dominating) solution  $\mathbf{a}^*$  such  
 121 that  $\text{visited}(\mathbf{a}^*) = \text{visited}(\mathbf{a})$ ,  $\text{current}(\mathbf{a}^*) = \text{current}(\mathbf{a})$  and  $\text{cost}(\mathbf{a}^*) < \text{cost}(\mathbf{a})$ . The tuple  
 122  $(\text{visited}(\mathbf{a}), \text{current}(\mathbf{a}))$  we refer to as the *DP state*, so removing all dominated partial solutions,  
 123 we keep exactly one minimum-cost solution for each unique DP state<sup>1</sup>. Since a solution can only  
 124 dominate other solutions with the same set of visited nodes, we only need to remove dominated  
 125 solutions from sets of solutions with the same number of actions. Therefore, we can efficiently  
 126 execute the DP algorithm in iterations, where at step  $t$  all solutions have (after  $t$  actions)  $t + 1$  visited  
 127 nodes (including the start node), keeping the memory need at  $O(B)$  states (with  $B$  the beam size).

128 We define the **scoring policy** using a pretrained model [26], which takes as input node coordinates  
 129 and edge distances to predict a raw ‘heatmap’ value  $\hat{h}_{ij} \in (0, 1)$  for each edge  $(i, j)$ . The model was  
 130 trained to predict optimal solutions, so  $\hat{h}_{ij}$  can be seen as the probability that edge  $(i, j)$  is in the  
 131 optimal tour. We force the heatmap to be symmetric thus we define  $h_{ij} = \max\{\hat{h}_{ij}, \hat{h}_{ji}\}$ . The policy  
 132 is defined using the heatmap values, in such a way to select the (partial) solutions with the largest  
 133 total *heat*, while also taking into account the (heat) *potential* for the unvisited nodes. The policy thus  
 134 selects the  $B$  solutions which have the highest *score*, defined as  $\text{score}(\mathbf{a}) = \text{heat}(\mathbf{a}) + \text{potential}(\mathbf{a})$ ,  
 135 with  $\text{heat}(\mathbf{a}) = \sum_{i=1}^{t-1} h_{a_{i-1}, a_i}$ , i.e. the sum of the heat of the edges, which can be computed  
 136 incrementally when expanding a solution. The potential is added as an estimate of the ‘heat-to-  
 137 go’ (similar to the heuristic in  $A^*$  search) for the remaining nodes, and avoids the ‘greedy pitfall’  
 138 of selecting the best edges while skipping over nearby nodes, which would prevent good edges  
 139 from being used later. It is defined as  $\text{potential}(\mathbf{a}) = \text{potential}_0(\mathbf{a}) + \sum_{i \notin \text{visited}(\mathbf{a})} \text{potential}_i(\mathbf{a})$   
 140 with  $\text{potential}_i(\mathbf{a}) = w_i \sum_{j \notin \text{visited}(\mathbf{a})} \frac{h_{ji}}{\sum_{k=0}^{n-1} h_{ki}}$ , where  $w_i$  is the node *potential weight* given by  
 141  $w_i = (\max_j h_{ji}) \cdot (1 - 0.1(\frac{c_{i0}}{\max_j c_{j0}} - 0.5))$ . By normalizing the heatmap values for incoming  
 142 edges, the (remaining) potential for node  $i$  is initially equal to  $w_i$  but decreases as good edges  
 143 become infeasible due to neighbours being visited. The node potential weight  $w_i$  is equal to the  
 144 maximum incoming edge heatmap value (an upper bound to the heat contributed by node  $i$ ), which  
 145 gets multiplied by a factor 0.95 to give a higher weight to nodes closer to the start node, which  
 146 we found helps to encourage the algorithm to keep edges that enable to return to the start node. The  
 147 overall heat + potential function identifies promising partial solutions and is computationally cheap.

### 148 3.2 Vehicle Routing Problem

149 For the VRP, we add a special depot node to the graph, indicated by DEP. Each node  $i$  has a demand  
 150  $d_i$ , and the goal is to find multiple routes, which have a limited capacity denoted by CAPACITY.

151 Additionally to the TSP **state variables**  $\text{cost}(\mathbf{a})$ ,  $\text{current}(\mathbf{a})$  and  $\text{visited}(\mathbf{a})$ , we keep track of  
 152  $\text{capacity}(\mathbf{a})$ , which is the *remaining* capacity in the current route/vehicle. A solution starts at the  
 153 depot, so we initialize the beam at step  $t = 0$  with the empty **initial solution** with  $\text{cost}(\mathbf{a}) = 0$ ,  
 154  $\text{current}(\mathbf{a}) = \text{DEP}$ ,  $\text{visited}(\mathbf{a}) = \emptyset$  and  $\text{capacity}(\mathbf{a}) = \text{CAPACITY}$ . For the VRP, we do not consider  
 155 visiting the depot as a separate action. Instead, we define  $2n$  actions, where  $a_t \in \{0, \dots, 2n - 1\}$ .  
 156 The actions  $0, \dots, n - 1$  indicate a *direct* move from the current node to node  $a_t$ , whereas the actions

<sup>1</sup>If we have multiple partial solutions with the same state and cost, we can arbitrarily choose one to dominate the other(s), for example the one with the lowest index of the current node.

157  $n, \dots, 2n - 1$  indicate a move to node  $a_t - n$  via the depot. **Feasible actions** are those that move  
 158 to unvisited nodes via edges in the graph and obey the following constraints. For the first action  
 159  $a_0$  there is no choice and we constrain (for convenience of implementation)  $a_0 \in \{n, \dots, 2n - 1\}$ .  
 160 A direct move ( $a_t < n$ ) is only feasible if  $d_{a_t} \leq \text{capacity}(\mathbf{a})$  and updates the state similar to TSP  
 161 but reduces remaining capacity by  $d_{a_t}$ . A move via the depot is always feasible (respecting the  
 162 graph edges and assuming  $d_i \leq \text{CAPACITY} \forall i$ ) as it resets the vehicle CAPACITY before subtracting  
 163 demand, but incurs the ‘via-depot cost’  $c_{ij}^{\text{DEP}} = c_{i,\text{DEP}} + c_{\text{DEP},j}$ . When all nodes are visited, we allow a  
 164 special action to return to the depot. This somewhat unusual way of representing a CVRP solution  
 165 has desirable properties similar to the TSP formulation: at step  $t$  we have exactly  $t$  nodes visited, and  
 166 we can run the DP in iterations, removing dominated solutions at each step  $t$ .

167 For VRP, a partial solution  $\mathbf{a}$  is a **dominated solution** dominated by  $\mathbf{a}^*$  if  $\text{visited}(\mathbf{a}^*) = \text{visited}(\mathbf{a})$   
 168 and  $\text{current}(\mathbf{a}^*) = \text{current}(\mathbf{a})$  (i.e.  $\mathbf{a}^*$  corresponds to the same DP state) and  $\text{cost}(\mathbf{a}^*) \leq \text{cost}(\mathbf{a})$   
 169 and  $\text{capacity}(\mathbf{a}^*) \geq \text{capacity}(\mathbf{a})$ , with *at least one of the two inequalities being strict*. This means  
 170 that for each DP state, given by the set of visited nodes and the current node, we do not only keep  
 171 the (single) solution with lowest cost (as in the TSP algorithm), but keep the complete set of pareto-  
 172 efficient solutions in terms of cost and remaining vehicle capacity. This is because a higher cost  
 173 partial solution may still be preferred if it has more remaining vehicle capacity, and vice versa.

174 For the VRP **scoring policy**, we modify the model [26] to include the depot node and demands. The  
 175 special depot node gets a separate learned initial embedding parameter, and we add additional edge  
 176 types for connections to the depot, to mark the depot as being special. Additionally, each node gets  
 177 an extra input (next to its coordinates) corresponding to  $d_i/\text{CAPACITY}$  (where we set  $d_{\text{DEP}} = 0$ ).  
 178 Apart from this, the model remains exactly the same<sup>2</sup>. The model is trained on example solutions  
 179 from LKH [24] (see Section 4.2), which are not optimal, but still provide a useful training signal.  
 180 Compared to TSP, the definition of the heat is slightly changed to accommodate for the ‘via-depot  
 181 actions’ and is best defined incrementally using the ‘via-depot heat’  $h_{ij}^{\text{DEP}} = h_{i,\text{DEP}} \cdot h_{\text{DEP},j} \cdot 0.1$ ,  
 182 where multiplication is used to keep heat values interpretable as probabilities and in the range  $(0, 1)$ .  
 183 The additional penalty factor of 0.1 for visiting the depot encourages the algorithm to minimize the  
 184 number of vehicles/routes. The initial heat is 0 and when expanding a solution  $\mathbf{a}$  to  $\mathbf{a}'$  using action  
 185  $a_t$ , the heat is incremented with either  $h_{\text{current}(\mathbf{a}),a_t}$  (if  $a_t < n$ ) or  $h_{\text{current}(\mathbf{a}),a_t-n}^{\text{DEP}}$  (if  $a_t \geq n$ ). The  
 186 potential is defined similarly to TSP, replacing the start node 0 by DEP.

### 187 3.3 Travelling Salesman Problem with Time Windows

188 For the TSPTW, we also have a special depot / start node 0, and each node  $i$  has a time window  
 189 defined by  $(l_i, u_i)$  in which the node should be visited, assuming travel time is equal to cost/distance.  
 190 It is allowed to wait if arrival at a node is before  $l_i$ , but arrival cannot be after  $u_i$  (i.e. the constraint is  
 191 hard). We consider the objective to be minimizing total *cost*, but minimizing total time (or *makespan*)  
 192 only requires training on different example solutions. Due to the hard constraints, TSPTW is typically  
 193 considered more challenging to solve than plain TSP, for which every solution is feasible.

194 The **state variables** and **initial solution** are equal to TSP except that we add  $\text{time}(\mathbf{a})$  which is  
 195 initially 0 ( $= l_0$ ). **Feasible actions**  $a_t \in \{0, \dots, n - 1\}$  are those that move to unvisited nodes via  
 196 edges in the graph such that the arrival time is no later than  $u_{a_t}$  and do not directly eliminate the  
 197 possibility to visit other nodes in time<sup>3</sup>. Expanding a solution  $\mathbf{a}$  to  $\mathbf{a}'$  using action  $a_t$  updates the  
 198 time as  $\text{time}(\mathbf{a}') = \max\{\text{time}(\mathbf{a}) + c_{\text{current}(\mathbf{a}),a_t}, l_{a_t}\}$ .

199 For each DP state, we keep all efficient solutions in terms of cost and time, so a partial solution  $\mathbf{a}$  is a  
 200 **dominated solution** dominated by  $\mathbf{a}^*$  if  $\mathbf{a}^*$  has the same DP state ( $\text{visited}(\mathbf{a}^*) = \text{visited}(\mathbf{a})$  and  
 201  $\text{current}(\mathbf{a}^*) = \text{current}(\mathbf{a})$ ) and is strictly better in terms of cost and time, i.e.  $\text{cost}(\mathbf{a}^*) \leq \text{cost}(\mathbf{a})$   
 202 and  $\text{time}(\mathbf{a}^*) \leq \text{time}(\mathbf{a})$ , with *at least one of the two inequalities being strict*.

203 The model [26] for the **scoring policy** is adapted to include the time windows  $(l_i, u_i)$  as node features  
 204 (in the same unit as coordinates and costs), and we use a special embedding for the depot similar to  
 205 VRP. Due to the time dimension, a TSPTW solution is *directed*, and edge  $(i, j)$  may be good whereas  
 206  $(j, i)$  may be not, so we adapt the model to enable predictions  $h_{ij} \neq h_{ji}$  (see details in Appendix  
 207 B). We generated example training solutions using (heuristic) DP with a large beam size, which was  
 208 faster than using LKH. Given the heat predictions, the score (heat + potential) is exactly as for TSP.

<sup>2</sup>Except that we do not use the K-nearest neighbour feature [26] as it contains no additional information.

<sup>3</sup>E.g., arriving at a node  $i$  at  $t = 10$  (including waiting) is not feasible if node  $j$  has  $u_j = 12$  and  $c_{ij} = 3$ .

209 **3.4 Graph sparsity**

210 As described, the DP algorithm can take into account a sparse graph to define feasible expansions. As  
 211 our problems are defined on sets of nodes rather than graphs, the use of a sparse graph is an artificial  
 212 design choice, which allows to significantly reduce the runtime but may sacrifice the possibility to  
 213 find good or optimal tours. We propose two different strategies for defining the sparse graph on  
 214 which to run the DP: thresholding the heatmap values  $h_{ij}$  and using the K-nearest neighbour (KNN)  
 215 graph. By default, we use a (low) heatmap threshold of  $10^{-5}$ , which rules out most of the edges as  
 216 the model confidently predicts (close to) 0 for most edges. This is a secondary way to leverage the  
 217 neural network (independent of the scoring policy), which can be seen as a form of learned *problem*  
 218 *reduction* [49]. For symmetric problems (TSP and VRP), we add KNN edges in both directions. For  
 219 the VRP, we additionally connect each node to the depot (and vice versa) to ensure feasibility.

220 **3.5 Implementation & hyperparameters**

221 We implement DPDP using PyTorch [44] to leverage batched computation on the GPU. For details,  
 222 see Appendix A. Our code is publicly available.<sup>4</sup> DPDP has very few hyperparameters, but the  
 223 heatmap threshold of  $10^{-5}$  and details like the functional form of e.g. the scoring policy are ‘educated  
 224 guesses’ or manually tuned on a few validation instances and can likely be improved. The runtime is  
 225 influenced by implementation choices which were manually selected using a few validation instances.

226 **4 Experiments**

227 **4.1 Travelling Salesman Problem**

228 In Table 1 we report our main results for DPDP with beam sizes of 10K (10 thousand) and 100K,  
 229 for the TSP with 100 nodes on a commonly used test set [29]. We report results using Concorde [2],  
 230 LKH [24] and Gurobi [22], as well as recent results of the strongest methods using neural networks  
 231 (‘neural approaches’) from literature. Running times for solving 10000 instances *after training* should  
 232 be taken as rough indications as some are on different machines, typically with 1 GPU or a many-core  
 233 CPU (8 - 32). The costs indicated with \* are not directly comparable due to slight dataset differences  
 234 [17]. Times for generating heatmaps (if applicable) is reported separately (as the first term) from the  
 235 running time for MCTS [17] or DP. DPDP achieves close to optimal results, strictly outperforming  
 236 the neural baselines achieving better results in less time (except POMO [30], see Section 4.2).

<sup>4</sup><https://github.com/?????>, to be disclosed after review

Table 1: Mean cost, gap and *total time* to solve 10000 TSP/VRP instances *after training*.

PROBLEM METHOD	TSP100			VRP100		
	COST	GAP	TIME	COST	GAP	TIME
CONCORDE [2]	7.765	0.000 %	6M			
HYBRID GENETIC SEARCH [56, 55]				15.563	0.000 %	6H11M
GUROBI [22]	7.776	0.151 %	31M			
LKH [24]	7.765	0.000 %	42M	15.647	0.536 %	12H57M
GNN HEATMAP + BEAM SEARCH [26]	7.87	1.39 %	40M			
LEARNING 2-OPT HEURISTICS [9]	7.83	0.87 %	41M			
MERGED GNN HEATMAP + MCTS [17]	7.764*	0.04 %	4M + 11M			
ATTENTION MODEL + SAMPLING [29]	7.94	2.26 %	1H	16.23	4.28 %	2H
STEP-WISE ATTENTION MODEL [60]	8.01	3.20 %	29S	16.49	5.96 %	39S
LEARNING IMPROV. HEURISTICS [59]	7.87	1.42 %	2H	16.03	3.00 %	5H
ATTENTION MODEL + POMO [30]	7.77	0.14 %	1M	15.76	1.26 %	2M
NEUREWRITER [7]				16.10	3.45 %	1H
DYNAMIC ATTN. MODEL + 2-OPT [45]				16.27	4.54 %	6H
NEUR. LRG. NEIGHB. SEARCH [25]				15.99	2.74 %	1H
LEARN TO IMPROVE [35]				15.57*	-	4000H
DPDP 10K	7.765	0.009 %	10M + 16M	15.830	1.713 %	10M + 50M
DPDP 100K	7.765	0.004 %	10M + 2H35M	15.694	0.843 %	10M + 5H48M
DPDP 1M				15.627	0.409 %	10M + 48H27M

Table 2: Mean cost, gap and *total time* to solve 10000 realistic [51] VRP100 instances after training.

METHOD	COST	GAP	TIME (1 GPU OR 16 CPUS)	TIME (4 GPUS OR 32 CPUS)
HGS [56, 55]	18050	0.000 %	7H53M	3H56M
LKH [24]	18133	0.507 %	25H32M	12H46M
DPDP 10K	18414	2.018 %	10M + 50M	2M + 13M
DPDP 100K	18253	1.127 %	10M + 5H48M	2M + 1H27M
DPDP 1M	18168	0.659 %	10M + 48H27M	2M + 12H7M

## 237 4.2 Vehicle Routing Problem

238 For the VRP, we train the model using 1 million instances of 100 nodes, generated according to the  
 239 distribution described by [41] and solved using one run of LKH [24]. We train using a batch size of  
 240 48 and a learning rate of  $10^{-3}$  (selected as the result of manual trials to best use our GPUs), for (at  
 241 most) 1500 epochs of 500 training steps (following [26]) from which we select the saved checkpoint  
 242 with the lowest validation loss. We use the validation and test sets by [29].

243 Table 1 shows the results compared to a recent implementation of Hybrid Genetic Search (HGS)<sup>5</sup>,  
 244 a SOTA heuristic VRP solver [56, 55]. HGS is faster and improves around 0.5% over LKH, which  
 245 is typically considered the baseline in related work. We present the results for LKH, as well as the  
 246 strongest neural approaches and DPDP with beam sizes up to 1 million. Some results used 2000  
 247 (different) instances [35] and cannot be directly compared<sup>6</sup>. DPDP outperforms all other neural  
 248 baselines, except POMO [30], which delivers good results very quickly by exploiting symmetries in  
 249 the problem. However, as it cannot (easily) improve further with additional runtime, we consider this  
 250 contribution orthogonal to DPDP. DPDP is competitive to LKH (see also Section 4.4).

251 **More realistic instances** We also train the model and run experiments with instances with 100  
 252 nodes from a more realistic and challenging data distribution [51]. This distribution, commonly used  
 253 in the routing community, has greater variability, in terms of node clustering and demand distributions.  
 254 LKH failed to solve two of the test instances, which we found out is because LKH by default uses  
 255 a fixed number of routes equal to a lower bound, given by  $\left\lceil \frac{\sum_{i=0}^{n-1} d_i}{\text{CAPACITY}} \right\rceil$ , which may be infeasible<sup>7</sup>.  
 256 Therefore we solve these instances by rerunning LKH with an unlimited number of allowed routes  
 257 (which in general gives worse results, see Section 4.4).

258 DPDP was run on a machine with 4 GPUs, but we also report (estimated) runtimes for 1 GPU  
 259 (1080Ti), and we compare against 16 or 32 CPUs for HGS and LKH. In Table 2 it can be seen that  
 260 the difference with LKH is, as expected, slightly larger than for the simpler dataset, but still below  
 261 1% for beam sizes of 100K-1M. We also observed a higher validation loss, so it may be possible to  
 262 improve results using more training data. Nevertheless, finding solutions within 1% of the specialized  
 263 SOTA HGS algorithm, and even closer to LKH, is impressive for these challenging instances, and we  
 264 consider the runtime (for solving 10K instances) acceptable, especially when using multiple GPUs.

## 265 4.3 TSP with Time Windows

266 For the TSP with hard time window constraints, we use the data distribution by [6] and use their set  
 267 of 100 test instances with 100 nodes. These were generated with small time windows, resulting in  
 268 a small feasible search space, such that even with very small beam sizes, our DP implementation  
 269 solves these instances optimally, eliminating the need for a policy. Therefore, we also consider a  
 270 more difficult distribution similar to [10], which has larger time windows which are more difficult  
 271 as the feasible search space is larger<sup>8</sup> [15]. For details, see Appendix B. For both distributions, we  
 272 generate training data and train the model exactly as we did for the VRP.

<sup>5</sup><https://github.com/vidalt/HGS-CVRP>

<sup>6</sup>The running time of 4000 hours (167 days) for 10K instances is estimated from 24min/instance [35].

<sup>7</sup>For example, three nodes with a demand of two cannot be assigned to two routes with a capacity of three.

<sup>8</sup>Up to a limit, as making the time windows infinite size reduces the problem to plain TSP.

Table 3: Mean cost, gap and *total time* to solve TSPTW100 instances after training.

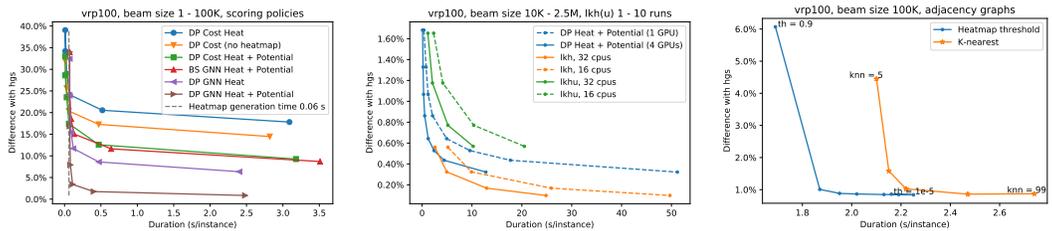
PROBLEM METHOD	SMALL TIME WINDOWS [6] (100 INST.)				LARGE TIME WINDOWS [10] (10K INST.)			
	COST	GAP	FAIL	TIME	COST	GAP	FAIL	TIME
GVNS 30x [10]	5129.58	0.000 %		7s	2432.112	0.000 %		37M15s
GVNS 1x [10]	5129.58	0.000 %		<1s	2457.974	1.063 %		1M4s
LKH 1x [24]	5130.32	0.014 %	1.00 %	5M48s	2431.404	-0.029 %		34H58M
BaB-DQN* [6]	5130.51	0.018 %		25H				
ILDS-DQN* [6]	5130.45	0.017 %		25H				
DPDP 10K	5129.58	0.000 %		6s + 1s	2431.143	-0.040 %		10M + 8M7s
DPDP 100K	5129.58	0.000 %		6s + 1s	2430.880	-0.051 %		10M + 1H16M

273 Table 3 shows the results for both data distributions, which are reported in terms of the difference  
 274 to General Variable Neighbourhood Search (GVNS) [10], the best open-source solver for TSPTW  
 275 we could find<sup>9</sup>, using 30 runs. For the small time window setting, both GVNS and DPDP find  
 276 optimal solutions for all 100 instances in just 7 seconds (in total, either on 16 CPUs or a single GPU).  
 277 LKH fails to solve one instance, but finds close to optimal solutions, but around 50 times slower.  
 278 BaB-DQN\* and ILDS-DQN\* [6], methods combining an existing solver with an RL trained neural  
 279 policy, take around 15 minutes *per instance* (orders of magnitudes slower) to solve most instances to  
 280 optimality. Due to complex set-up, we were unable to run BaB-DQN\* and ILDS-DQN\* ourselves for  
 281 the setting with larger time windows. In this setting, we find DPDP outperforms both LKH (where  
 282 DPDP is orders of magnitude faster) and GVNS, in both speed and solution quality. This illustrates  
 283 that DPDP, due to its nature, is especially well suited to handle constrained problems.

#### 284 4.4 Ablations

285 **Scoring policy** To evaluate the value of different components of DPDP’s **GNN Heat + Potential**  
 286 scoring policy, we compare against other variants. **GNN Heat** is the version without the potential,  
 287 whereas **Cost Heat + Potential** and **Cost Heat** are variants that use a ‘heuristic’  $\hat{h}_{ij} = \frac{c_{ij}}{\max_k c_{ik}}$   
 288 instead of the GNN. **Cost** directly uses the current cost of the solution, and can be seen as ‘classic’  
 289 restricted DP. Finally, **BS GNN Heat + Potential** uses beam search without dynamic programming,  
 290 i.e. without removing dominated solutions. To evaluate only the scoring policy, each variant uses  
 291 the fully connected graph ( $knn = n - 1$ ). Figure 3a shows the value of DPDP’s potential function,  
 292 although even without it results are still significantly better than ‘classic’ heuristic DP variants using  
 293 cost-based scoring policies. Also, it is clear that using DP significantly improves over a standard beam  
 294 search (by removing dominated solutions). Lastly, the figure illustrates how the time for generating  
 295 the heatmap using the neural network, despite its significant value, only makes up a small portion of  
 296 the total runtime.

<sup>9</sup><https://github.com/sashakh/TSPTW>



(a) Different scoring policies, as well as ‘pure’ beam search, for beam sizes 1, 10, 100, 1000, 10K, 100K. (b) Beam sizes 10K, 25K, 50K, 100K, 250K, 500K, 1M, 2.5M compared against LKH(U) with 1, 2, 5 and 10 runs. (c) Sparsities with heatmap thresholds 0.9, 0.5, 0.2, 0.1,  $10^{-2}$ ,  $10^{-3}$ ,  $10^{-4}$ ,  $10^{-5}$  and  $knn = 5, 10, 20, 50, 99$ . Beam size 100K.

Figure 3: DPDP ablations on 100 validation instances of VRP with 100 nodes.

297 **Beam size** DPDP allows to trade off the performance vs. the runtime using the beam size  $B$  (and to  
 298 some extent the graph sparsity, see Section 4.4). We illustrate this trade-off in Figure 3b, where we  
 299 evaluate DPDP on 100 validation instances for VRP, with different beam sizes from 10K to 2.5M.  
 300 We also report the trade-off curve for the LKH(U), which is the strongest baseline that can also  
 301 solve different problems. We vary the runtime using 1, 2, 5 and 10 runs (returning the best solution).  
 302 LKHU(nlimited) is the version which allows an unlimited number of routes (see Section 4.2). It is  
 303 hard to compare GPU vs CPU, so we report (estimated) runtimes for different hardware, i.e. 1 or 4  
 304 GPUs (with 3 CPUs per GPU) and 16 or 32 CPUs. We report the difference (i.e. the gap) with HGS,  
 305 analog to how results are reported in Table 1. We emphasize that in most related work (e.g. [29]), the  
 306 strongest baseline considered is one run of LKH, so we compare against a much stronger baseline.  
 307 Also, our goal is not to outperform HGS (which is SOTA and specific to VRP) or LKH, but to show  
 308 DPDP has reasonable performance, while being a flexible framework for other (routing) problems.

309 **Graph sparsity** We test the two graph sparsification strategies described in Section 3.4 as another  
 310 way to trade off performance and runtime of DPDP. In Figure 3c, we experiment with different  
 311 heatmap thresholds from  $10^{-5}$  to 0.9 and different values for KNN from 5 to 99 (fully connected).  
 312 The heatmap threshold strategy clearly outperforms the KNN strategy as it yields the same results  
 313 using sparser graphs (and lower runtimes). This illustrates that the heatmap threshold strategy is more  
 314 informed than the KNN strategy, confirming the value of the neural network predictions.

## 315 5 Discussion

316 In this paper we introduced Deep Policy Dynamic Programming, which combines machine learning  
 317 and dynamic programming for solving vehicle routing problems. The method yields close to optimal  
 318 results for TSPs with 100 nodes and is competitive to the highly optimized LKH [24] solver for VRPs  
 319 with 100 nodes. On the TSP with time windows, DPDP also outperforms LKH, being significantly  
 320 faster, as well as GVNS [10], the best open source solver we could find. Given that DPDP was not  
 321 specifically designed for TSPTW, and still has possibilities for improvement, we consider this an  
 322 impressive and promising achievement.

323 The constructive nature of DPDP (combined with search) allows to efficiently address hard constraints  
 324 such as time windows, which are typically considered challenging in neural combinatorial optimiza-  
 325 tion [5, 29] and are also difficult for local search heuristics (as they need to maintain feasibility while  
 326 adapting a solution). Given our results on TSP, VRP and TSPTW, and the flexibility of DP as a  
 327 framework, we think DPDP has great potential for solving many more variants of routing problems,  
 328 and possibly even other problems that can be formulated using DP (e.g. job shop scheduling [21]).  
 329 We hope that our work brings machine learning research for combinatorial optimization closer to  
 330 the operations research (especially vehicle routing) community, by combining machine learning  
 331 with DP and evaluating the resulting new framework on different data distributions used by different  
 332 communities [41, 51, 6, 10].

333 **Scope, limitations & future work** Deep learning for combinatorial optimization is a recent re-  
 334 search direction, which could significantly impact the way practical optimization problems get solved  
 335 in the future. Currently, however, it is still hard to beat most SOTA problem specific solvers from the  
 336 OR community. Despite our success for TSPTW, DPDP is not yet a practical alternative in general,  
 337 but we do consider our results as highly encouraging for further research. We believe such research  
 338 could yield significant further improvement by addressing key current limitations: (1) the scalability  
 339 to larger instances, (2) the dependency on example solutions and (3) the heuristic nature of the scoring  
 340 function. First, while 100 nodes is not far from the size of common benchmarks (100 - 1000 for VRP  
 341 [51] and 20 - 200 for TSPTW [10]), scaling is a challenge, mainly due to the ‘fully-connected’  $O(n^2)$   
 342 graph neural network. Future work could reduce this complexity following e.g. [33]. The dependency  
 343 on example solutions from an existing solver also becomes more prominent for larger instances,  
 344 but could potentially be removed by ‘bootstrapping’ using DP itself as we, in some sense, have  
 345 done for TSPTW (see Section 3.3). Future work could iterate this process to train the model ‘tabula  
 346 rasa’ (without example solutions), where DP could be seen analogous to MCTS in *AlphaZero* [48].  
 347 Lastly, the heat + potential score function is a well-motivated but heuristic function that was manually  
 348 designed as a function of the predicted heatmap. While it worked well for the three problems we  
 349 considered, it may need suitable adaption for other problems. Training this function end-to-end  
 350 [11, 58], while keeping a low computational footprint, would be an interesting topic for future work.

## References

- [1] Luca Accorsi and Daniele Vigo. A fast and scalable heuristic for the solution of large-scale capacitated vehicle routing problems. Technical report, Tech. rep., University of Bologna, 2020.
- [2] David Applegate, Robert Bixby, Vasek Chvatal, and William Cook. Concorde TSP solver, 2006.
- [3] Richard Bellman. On the theory of dynamic programming. *Proceedings of the National Academy of Sciences of the United States of America*, 38(8):716, 1952.
- [4] Richard Bellman. Dynamic programming treatment of the travelling salesman problem. *Journal of the ACM (JACM)*, 9(1):61–63, 1962.
- [5] Irwan Bello, Hieu Pham, Quoc V Le, Mohammad Norouzi, and Samy Bengio. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.
- [6] Quentin Cappart, Thierry Moisan, Louis-Martin Rousseau, Isabelle Prémont-Schwarz, and Andre Cire. Combining reinforcement learning and constraint programming for combinatorial optimization. *arXiv preprint arXiv:2006.01610*, 2020.
- [7] Xinyun Chen and Yuandong Tian. Learning to perform local rewriting for combinatorial optimization. In *Advances in Neural Information Processing Systems*, volume 32, pages 6281–6292, 2019.
- [8] William Cook and Paul Seymour. Tour merging via branch-decomposition. *INFORMS Journal on Computing*, 15(3):233–248, 2003.
- [9] Paulo Roberto de O da Costa, Jason Rhuggenaath, Yingqian Zhang, and Alp Akcay. Learning 2-opt heuristics for the traveling salesman problem via deep reinforcement learning. *Proceedings of Machine Learning Research*, 1:17, 2020.
- [10] Rodrigo Ferreira Da Silva and Sebastián Urrutia. A general vns heuristic for the traveling salesman problem with time windows. *Discrete Optimization*, 7(4):203–211, 2010.
- [11] Hal Daumé III and Daniel Marcu. Learning as search optimization: Approximate large margin methods for structured prediction. In *Proceedings of the 22nd international conference on Machine learning*, pages 169–176, 2005.
- [12] Arthur Delarue, Ross Anderson, and Christian Tjandraatmadja. Reinforcement learning with combinatorial actions: An application to vehicle routing. *Advances in Neural Information Processing Systems*, 33, 2020.
- [13] Michel Deudon, Pierre Cournut, Alexandre Lacoste, Yossiri Adulyasak, and Louis-Martin Rousseau. Learning heuristics for the TSP by policy gradient. In *International Conference on the Integration of Constraint Programming, Artificial Intelligence, and Operations Research*, pages 170–181. Springer, 2018.
- [14] Edsger W Dijkstra. A note on two problems in connexion with graphs. *Numerische mathematik*, 1(1):269–271, 1959.
- [15] Yvan Dumas, Jacques Desrosiers, Eric Gelinas, and Marius M Solomon. An optimal algorithm for the traveling salesman problem with time windows. *Operations research*, 43(2):367–371, 1995.
- [16] Jonas K Falkner and Lars Schmidt-Thieme. Learning to solve vehicle routing problems with time windows through joint attention. *arXiv preprint arXiv:2006.09100*, 2020.
- [17] Zhang-Hua Fu, Kai-Bin Qiu, and Hongyuan Zha. Generalize a small pre-trained model to arbitrarily large tsp instances. *arXiv preprint arXiv:2012.10658*, 2020.
- [18] Lei Gao, Mingxiang Chen, Qichang Chen, Ganzhong Luo, Nuoyi Zhu, and Zhixin Liu. Learn to design the heuristics for vehicle routing problem. *arXiv preprint arXiv:2002.08539*, 2020.
- [19] Maxime Gasse, Didier Chetelat, Nicola Ferroni, Laurent Charlin, and Andrea Lodi. Exact combinatorial optimization with graph convolutional neural networks. In *Advances in Neural Information Processing Systems*.

- 398 [20] Joaquim Gromicho, Jelke J van Hoorn, Adrianus Leendert Kok, and Johannes MJ Schutten.  
399 Restricted dynamic programming: a flexible framework for solving realistic vrps. *Computers &*  
400 *operations research*, 39(5):902–909, 2012.
- 401 [21] Joaquim AS Gromicho, Jelke J Van Hoorn, Francisco Saldanha-da Gama, and Gerrit T Timmer.  
402 Solving the job-shop scheduling problem optimally by dynamic programming. *Computers &*  
403 *Operations Research*, 39(12):2968–2977, 2012.
- 404 [22] Gurobi Optimization, LLC. Gurobi, 2018.
- 405 [23] Michael Held and Richard M Karp. A dynamic programming approach to sequencing problems.  
406 *Journal of the Society for Industrial and Applied Mathematics*, 10(1):196–210, 1962.
- 407 [24] Keld Helsgaun. An extension of the Lin-Kernighan-Helsgaun TSP solver for constrained  
408 traveling salesman and vehicle routing problems: Technical report. 2017.
- 409 [25] André Hottung and Kevin Tierney. Neural large neighborhood search for the capacitated vehicle  
410 routing problem. *arXiv preprint arXiv:1911.09539*, 2019.
- 411 [26] Chaitanya K Joshi, Thomas Laurent, and Xavier Bresson. An efficient graph convolutional  
412 network technique for the travelling salesman problem. *arXiv preprint arXiv:1906.01227*, 2019.
- 413 [27] Chaitanya K Joshi, Thomas Laurent, and Xavier Bresson. On learning paradigms for the  
414 travelling salesman problem. *arXiv preprint arXiv:1910.07210*, 2019.
- 415 [28] AL Kok, Elias W Hans, Johannes MJ Schutten, and Willem HM Zijm. A dynamic programming  
416 heuristic for vehicle routing with time-dependent travel times and required breaks. *Flexible*  
417 *services and manufacturing journal*, 22(1-2):83–108, 2010.
- 418 [29] Wouter Kool, Herke van Hoof, and Max Welling. Attention, learn to solve routing problems! In  
419 *International Conference on Learning Representations*, 2019.
- 420 [30] Yeong-Dae Kwon, Jinho Choo, Byoungjip Kim, Iljoo Yoon, Youngjune Gwon, and Seungjai  
421 Min. Pomo: Policy optimization with multiple optima for reinforcement learning. *Advances in*  
422 *Neural Information Processing Systems*, 33, 2020.
- 423 [31] Gilbert Laporte. The vehicle routing problem: An overview of exact and approximate algorithms.  
424 *European journal of operational research*, 59(3):345–358, 1992.
- 425 [32] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436,  
426 2015.
- 427 [33] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh.  
428 Set transformer: A framework for attention-based permutation-invariant neural networks. In  
429 *International Conference on Machine Learning*, pages 3744–3753. PMLR, 2019.
- 430 [34] Zhuwen Li, Qifeng Chen, and Vladlen Koltun. Combinatorial optimization with graph convolu-  
431 tional networks and guided tree search. *Advances in Neural Information Processing Systems*,  
432 page 539, 2018.
- 433 [35] Hao Lu, Xingwen Zhang, and Shuang Yang. A learning-based iterative method for solving  
434 vehicle routing problems. In *International Conference on Learning Representations*, 2020.
- 435 [36] Qiang Ma, Suwen Ge, Danyang He, Darshan Thaker, and Iddo Drori. Combinatorial opti-  
436 mization by graph pointer networks and hierarchical reinforcement learning. *arXiv preprint*  
437 *arXiv:1911.04936*, 2019.
- 438 [37] Chryssi Malandraki and Robert B Dial. A restricted dynamic programming heuristic algorithm  
439 for the time dependent traveling salesman problem. *European Journal of Operational Research*,  
440 90(1):45–55, 1996.
- 441 [38] Nina Mazyavkina, Sergey Sviridov, Sergei Ivanov, and Evgeny Burnaev. Reinforcement learning  
442 for combinatorial optimization: A survey. *arXiv preprint arXiv:2003.03600*, 2020.

- 443 [39] Aristide Mingozzi, Lucio Bianco, and Salvatore Ricciardelli. Dynamic programming strategies  
444 for the traveling salesman problem with time window and precedence constraints. *Operations*  
445 *research*, 45(3):365–377, 1997.
- 446 [40] Vinod Nair, Sergey Bartunov, Felix Gimeno, Ingrid von Glehn, Pawel Lichocki, Ivan Lobov,  
447 Brendan O’Donoghue, Nicolas Sonnerat, Christian Tjandraatmadja, Pengming Wang, et al.  
448 Solving mixed integer programs using neural networks. *arXiv preprint arXiv:2012.13349*, 2020.
- 449 [41] MohammadReza Nazari, Afshin Oroojlooy, Lawrence Snyder, and Martin Takac. Reinforcement  
450 learning for solving the vehicle routing problem. In *Advances in Neural Information Processing*  
451 *Systems*, pages 9860–9870, 2018.
- 452 [42] Clara Novoa and Robert Storer. An approximate dynamic programming approach for the  
453 vehicle routing problem with stochastic demands. *European Journal of Operational Research*,  
454 196(2):509–515, 2009.
- 455 [43] Alex Nowak, Soledad Villar, Afonso S Bandeira, and Joan Bruna. A note on learning algorithms  
456 for quadratic assignment with graph neural networks. *arXiv preprint arXiv:1706.07450*, 2017.
- 457 [44] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,  
458 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in  
459 pytorch. 2017.
- 460 [45] Bo Peng, Jiahai Wang, and Zizhen Zhang. A deep reinforcement learning algorithm using dy-  
461 namic attention model for vehicle routing problems. In *International Symposium on Intelligence*  
462 *Computation and Applications*, pages 636–650. Springer, 2019.
- 463 [46] Stefan Ropke and David Pisinger. An adaptive large neighborhood search heuristic for the  
464 pickup and delivery problem with time windows. *Transportation science*, 40(4):455–472, 2006.
- 465 [47] Gerhard Schrimpf, Johannes Schneider, Hermann Stamm-Wilbrandt, and Gunter Dueck. Record  
466 breaking optimization results using the ruin and recreate principle. *Journal of Computational*  
467 *Physics*, 159(2):139–171, 2000.
- 468 [48] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur  
469 Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. A general  
470 reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*,  
471 362(6419):1140–1144, 2018.
- 472 [49] Yuan Sun, Andreas Ernst, Xiaodong Li, and Jake Weiner. Generalization of machine learning  
473 for problem reduction: a case study on travelling salesman problems. *OR Spectrum*, pages 1–27,  
474 2020.
- 475 [50] Paolo Toth and Daniele Vigo. *Vehicle routing: problems, methods, and applications*. SIAM,  
476 2014.
- 477 [51] Eduardo Uchoa, Diego Pecin, Artur Pessoa, Marcus Poggi, Thibaut Vidal, and Anand Subrama-  
478 nian. New benchmark instances for the capacitated vehicle routing problem. *European Journal*  
479 *of Operational Research*, 257(3):845–858, 2017.
- 480 [52] Wouter van Heeswijk and Han La Poutré. Approximate dynamic programming with neural  
481 networks in linear discrete action spaces. *arXiv preprint arXiv:1902.09855*, 2019.
- 482 [53] Jelke J. van Hoorn. *Dynamic Programming for Routing and Scheduling*. PhD thesis, 2016.
- 483 [54] Natalia Vesselinova, Rebecca Steinert, Daniel F Perez-Ramirez, and Magnus Boman. Learning  
484 combinatorial optimization on graphs: A survey with applications to networking. *IEEE Access*,  
485 8:120388–120416, 2020.
- 486 [55] Thibaut Vidal. Hybrid genetic search for the cvrp: Open-source implementation and swap\*  
487 neighborhood. *arXiv preprint arXiv:2012.10384*, 2020.
- 488 [56] Thibaut Vidal, Teodor Gabriel Crainic, Michel Gendreau, Nadia Lahrichi, and Walter Rei. A  
489 hybrid genetic algorithm for multidepot and periodic vehicle routing problems. *Operations*  
490 *Research*, 60(3):611–624, 2012.

- 491 [57] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Advances in Neural*  
492 *Information Processing Systems*, pages 2692–2700, 2015.
- 493 [58] Sam Wiseman and Alexander M Rush. Sequence-to-sequence learning as beam-search opti-  
494 mization. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language*  
495 *Processing*, pages 1296–1306, 2016.
- 496 [59] Yaoxin Wu, Wen Song, Zhiguang Cao, Jie Zhang, and Andrew Lim. Learning improvement  
497 heuristics for solving routing problems. *arXiv preprint arXiv:1912.05784*, 2019.
- 498 [60] Liang Xin, Wen Song, Zhiguang Cao, and Jie Zhang. Step-wise deep learning models for  
499 solving routing problems. *IEEE Transactions on Industrial Informatics*, 2020.
- 500 [61] Shenghe Xu, Shivendra S Panwar, Murali Kodialam, and TV Lakshman. Deep neural network  
501 approximated dynamic programming for combinatorial optimization. In *Proceedings of the*  
502 *AAAI Conference on Artificial Intelligence*, volume 34, pages 1684–1691, 2020.
- 503 [62] Feidiao Yang, Tiancheng Jin, Tie-Yan Liu, Xiaoming Sun, and Jialin Zhang. Boosting dynamic  
504 programming with neural networks for solving np-hard problems. In *Asian Conference on*  
505 *Machine Learning*, pages 726–739. PMLR, 2018.

## 506 Checklist

- 507 1. For all authors...
- 508 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
509 contributions and scope? [Yes]
- 510 (b) Did you describe the limitations of your work? [Yes] Scaling and scoring policy, see  
511 discussion
- 512 (c) Did you discuss any potential negative societal impacts of your work? [No] We do not  
513 see any issues
- 514 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
515 them? [Yes]
- 516 2. If you are including theoretical results...
- 517 (a) Did you state the full set of assumptions of all theoretical results? [Yes] For claim that  
518 DPDP is asymptotically optimal, the DP description defines the set of assumptions and  
519 we note the beam size for which this holds.
- 520 (b) Did you include complete proofs of all theoretical results? [No] Optimality and  
521 complexity of full DP for TSP is already known [23, 4] (we do not claim this is a new  
522 result, we just mention it)
- 523 3. If you ran experiments...
- 524 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
525 mental results (either in the supplemental material or as a URL)? [Yes] Code will be  
526 available upon release
- 527 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
528 were chosen)? [Yes] We did a fair job to include all the information necessary to  
529 reproduce the experiments
- 530 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
531 ments multiple times)? [No] For training the model, we used a standard supervised  
532 learning setup and we found training to be stable across different problems and data  
533 distributions. Due to high computational costs, we did not do multiple runs, which  
534 we consider conservative as likely results can be improved by further experiments and  
535 training with multiple seeds. As for evaluation of the algorithms on random instances,  
536 we use a large test set of 10,000 instances. As we use the same instances for our own  
537 implemented baselines, the statistical significance in such a 'paired' setting is much  
538 higher than would be suggested by error bars. For results directly reported from other  
539 papers, the margin is large enough with 10000 instances or we avoid strong claims.  
540 With our code, we plan to release all solutions for all test instances for the methods we  
541 evaluated ourselves.
- 542 (d) Did you include the total amount of compute and the type of resources used (e.g.,  
543 type of GPUs, internal cluster, or cloud provider)? [Yes] While not explicitly, the  
544 total computational requirement for generating training data, which requires the most  
545 compute, can be derived from solving 1M instances and the time reported for solving  
546 10K instances in Table 1. This is a number of days (around 4) on a cluster with 10  
547 32-core machines, per problem.
- 548 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 549 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 550 (b) Did you mention the license of the assets? [No] This can be found in the original  
551 repositories
- 552 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
553 Will be provided with the released code as URL
- 554 (d) Did you discuss whether and how consent was obtained from people whose data you're  
555 using/curating? [No] Data is public or provided in personal communication, consent  
556 for publication was given
- 557 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
558 information or offensive content? [N/A]
- 559 5. If you used crowdsourcing or conducted research with human subjects...

- 560 (a) Did you include the full text of instructions given to participants and screenshots, if  
561 applicable? [N/A]
- 562 (b) Did you describe any potential participant risks, with links to Institutional Review  
563 Board (IRB) approvals, if applicable? [N/A]
- 564 (c) Did you include the estimated hourly wage paid to participants and the total amount  
565 spent on participant compensation? [N/A]