# DEEP LEARNING-BASED HEURISTIC CONSTRUCTION FOR ROUTING PROBLEMS WITH DYNAMIC ENCODER AND DUAL-CHANNEL DECODER ARCHITECTURE

## Anonymous authors

Paper under double-blind review

### ABSTRACT

The routing problem is a classic combinatorial optimization challenge. Constructing heuristics using deep learning models presents a promising approach for its resolution. In this paper, we propose a novel model with a dynamic encoder and dual-channel decoder (DEDD) architecture to learn construction heuristics for the routing problem. The dynamic encoder encodes the node features of the decomposed subproblems at each selection step, thereby obtaining more accurate node embeddings. The dual-channel decoder facilitates more diverse node selections at each step, increasing the probability of the model identifying optimal solutions. Additionally, we design an effective node selection strategy to assist the model in choosing nodes at each step. Experimental results on the Traveling Salesman Problem (TSP) and the Capacitated Vehicle Routing Problem (CVRP) with up to 1000 nodes demonstrate that the solutions generated by the DEDD model are nearly optimal, underscoring its efficacy.

024 025 026

006

008 009 010

011 012 013

014

015

016

017

018

019

021

# 1 INTRODUCTION

027 028 029

The routing problem (Veres & Moussa, 2019) is a type of combinatorial optimization (CO) problem prevalent in logistics, distribution, transportation, and other fields with significant practical 031 application value (Toth & Vigo, 2014). Its fundamental problems and classic variants include the Traveling Salesman Problem (TSP) and the Capacitated Vehicle Routing Problem (CVRP). The 033 routing problem is NP-hard, making it particularly challenging to solve (Li et al., 2022). Current 034 approaches for solving CO problems are mainly divided into four categories: exact approaches, approximation algorithms, heuristic algorithms and deep learning-based approaches. Exact ap-035 proaches, such as branch-and-bound (Fischetti et al., 1994; Lawler & Wood, 1966) and dynamic 036 programming (Bertsekas, 2000; Sniedovich, 2010), can theoretically find the global optimal so-037 lution but have exponential worst-case computational complexity. Thus, as the problem size increases, the computational effort required by these algorithms becomes impractical. Approximation algorithms (Williamson & Shmoys, 2011; Vazirani, 2001; Hochba, 1997) can obtain a theoretically 040 guaranteed solution with polynomial computational complexity, but the solution quality is often 041 suboptimal. Heuristic algorithms (Van Laarhoven et al., 1992; Wesley Barnes & LAGUNA, 1993; 042 Halim & Ismail, 2019) involve experts designing specific heuristic rules based on extensive domain 043 knowledge to search the solution space within an acceptable timeframe to find feasible solutions to 044 CO problems. Due to the limitations of these three types of approaches, numerous deep learning models have been proposed in recent years to address some of the issues presented in the aforementioned approaches. 046

Learning-based neural combinatorial optimization approaches can be divided into two categories:
learning improvement heuristics (Chen & Tian, 2019; Wu et al., 2021; Ma et al., 2021) and learning
construction heuristics (Vinyals et al., 2015; Joshi et al., 2019; Kool et al., 2022). Learning improvement heuristics use deep reinforcement learning algorithms to learn the rules of iterative search operators, iteratively searching for solutions based on the learned rules. Essentially, learning improvement heuristics are iterative search algorithms with good optimization effects. The approaches proposed by Chen & Tian (2019) and Lu et al. (2019) have achieved results that match or even surpass professional combinatorial optimization solvers such as LKH3 (Helsgaun, 2017), Google OR tools

(Perron & Furnon, 2019), Gurobi, and Concorde (Applegate et al., 2006). However, these methods still lag behind end-to-end approaches.

Learning construction heuristics primarily use deep neural networks with an encoder-decoder struc-057 ture. The encoder maps node information to embeddings, while the decoder, based on embeddings 058 and other information, provides the probability of each node being selected at each step. Finally, 059 based on these probabilities, a rule (e.g., greedily selecting the node with the highest probability 060 Luo et al. (2023) determines the node chosen at each step. After providing the problem instance as 061 input, the model repeatedly adds nodes to partial solutions until a complete solution is generated. 062 Since this approach inputs the problem instance directly into the deep neural network and outputs 063 the solution directly, it is also called an "end-to-end method." Compared to iterative heuristic algo-064 rithms, learning construction heuristics offer several advantages: First, they directly output the solution, resulting in faster solving speeds. Second, deep neural networks learn heuristic rules from data, 065 replacing manually designed rules by experts. Finally, traditional CO problem-solving algorithms 066 usually run on CPU, whereas deep neural networks run on GPU, enabling better parallelization and 067 faster inference on large batches of problem instances. However, the quality of solutions directly 068 output by deep neural networks often has room for improvement. Thus, various approaches (e.g., 069 beam search (Nazari et al., 2018)) are employed to refine the initially generated solutions for better results. 071

To explore more ways to solve CO problems, scholars have proposed many deep learning models with different structures, trained through supervised learning (Fu et al., 2021; Joshi et al., 2022; Hottung et al., 2021a) or reinforcement learning (Hottung & Tierney, 2020; Hottung et al., 2021b; d O Costa et al., 2020). However, existing models still have some limitations. First, they are mostly confined to solving small-scale problems and perform poorly on larger-scale problems. This is because most models are trained and tested on specific scales (e.g., TSP100), leading to suboptimal generalization.

Second, existing learning-based neural combinatorial optimization models typically learn only one construction strategy, limiting the diversity of decisions when selecting nodes at each step. The decoder provides a probability matrix for node selection at each step, which is derived from the rules learned and the node embeddings output by the encoder. The model ultimately selects the node to visit based on the probability matrix and certain selection rules. In other words, the diversity of node selection strategies at each construction step mainly comes from the decoder's probability matrix and the node selection strategy based on it, which is insufficient for constructing high-quality solutions. Different node choices at each step lead to different directions for exploring complete solutions, so enriching the diversity of node selection strategies at each step is crucial.

Additionally, current deep learning models encode the problem instance into global node embeddings and other information embeddings once at the beginning. In subsequent construction steps, the decoder uses the same node embeddings to build the solution. However, when solving a routing problem with state transitions, using node embeddings from the previous state to solve the current state problem may lead to suboptimal strategies. Also, encoding all node at once may cause the model to learn scale-related features, performing well on trained scales but failing to capture necessary relationships among nodes in untrained scales.

094 To address these limitations, we propose a deep learning model with a dynamic encoder and dual-095 channel decoder (DEDD) structure for learning construction heuristics for routing problems. Firstly, 096 to enhance model performance, the encoder in the DEDD model dynamically selects different nodes 097 at each step to form new subproblems as inputs, allowing the decoder to choose the next node based 098 on more accurate embeddings. Secondly, we propose a reasonable subproblem construction ap-099 proach. At each construction step, we form subproblems from the starting point of the complete instance, the visited node from the previous step, and the remaining available nodes, balancing 100 global and local information. Furthermore, since the dynamic encoder's input changes with con-101 struction steps, the DEDD model tends to learn scale-independent features, making it less sensitive 102 to instance sizes and better at generalizing across different problem scales. Lastly, to enhance the 103 effectiveness of node selection strategies, we propose a dual-channel decoder structure. The dual 104 channels independently provide node probability matrices at each step, with the model selecting one 105 channel's result to update the partial solution based on the probability matrix and selection strat-106 egy designed by us. Different channels share parameters in all but the final attention layer, thus 107 enhancing strategy effectiveness while minimizing computational overhead.

It is worth noting that the DEDD model is not intended to completely surpass existing highly optimized professional routing problem solvers but to explore designing deep learning models that autonomously learn stronger heuristic rules for solving CO problems. We apply the model to solving TSP and CVRP of various scales, and experimental results demonstrate that the proposed DEDD model achieves good performance within shorter inference times and effectively solves instances up to a scale of 1000 nodes

114 115

# 2 RELATED WORKS

116 117

In 1985, Hopfield & Tank (1985) introduce the Hopfield network for solving the TSP and other CO 118 problems, pioneering the use of neural networks for CO problem-solving. However, the Hopfield 119 network requires retraining for each new TSP instance. In 2015, Vinyals et al. (2015) introduce the 120 Pointer Network (PtrNet), the first to effectively employ deep neural networks for CO problems. The 121 PtrNet features an encoder-decoder structure, with both components composed of long short-term 122 memory (LSTM) networks. The model constructs complete solutions in an autoregressive man-123 ner through supervised learning. Given that labels are the optimal solutions of problem instances, 124 and acquiring these solutions is costly, supervised learning becomes impractical. Bello et al. (2016) 125 propose using the efficient reinforcement learning approach A3C to replace supervised learning, al-126 lowing the model to be trained on larger problem instances. Nazari et al. (2018) also use A3C to train 127 and optimize the PtrNet. Kool et al. (2022) and Deudon et al. (2018) employ the Transformer architecture for CO problems. Kool et al. (2022) propose an attention model with an encoder-decoder 128 structure. The encoder has three layers of attention, encoding all nodes simultaneously during train-129 ing and inference. The decoder has a single layer of attention, utilizing a pointer-like attention 130 mechanism for decoding. Since the proposal of the attention model (AM), numerous learning con-131 struction heuristics based on AM have been developed, but these methods can only perform well on 132 small-scale CO problems. 133

In addition to deep learning models with encoder-decoder structures, some researchers have pro-134 posed approaches based on graph neural networks. Khalil et al. (2017) are the first to use the Q-135 Learning (Watkins & Dayan, 1992) algorithm of deep reinforcement learning to train graph neural 136 networks for solving CO problems, achieving strong results in minimum vertex cover and maximum 137 cut, but their results for TSP are less ideal. Ma et al. (2019) combine the PtrNet and graph neural 138 networks by transforming graph neural network node features into node embeddings and using the 139 PtrNet's attention mechanism to construct complete solutions, achieving strong optimization results 140 in TSP. However, hierarchical reinforcement learning requires the setting of goals in order to achieve 141 good performance, and the model has low efficiency in exploration during training. 142

Some researchers have used deep reinforcement learning to enhance search algorithms. Traditional search algorithms typically involve experts manually designing heuristic rules based on specialized domain knowledge. Currently, some researchers use deep reinforcement learning to enable models to automatically learn or select heuristic rules, iteratively searching for solutions based on learned



Figure 1: The architecture of the proposed DEDD model, featuring a single-layer encoder and a dual-channel decoder.

162 rules and obtaining higher-quality solutions after multiple iterations. Chen & Tian (2019) propose 163 NeuRewriter to learn improvement heuristics, training two heuristic strategies to recursively re-164 fine initial solutions, achieving or surpassing the performance of professional solvers like LKH3 165 (Helsgaun, 2017). However, a limitation of this approach is the difficulty in parallelizing the solv-166 ing of instances. Yolcu & Póczos (2019) adopt a local search framework, using deep reinforcement learning to learn variable selection operators, finding better-quality solutions in fewer iterations but 167 with a long runtime. Additionally, many approaches that iteratively improve solutions still rely 168 on expert-designed heuristic rules, and the iterative steps cannot be parallelized. Thus, learning improvement heuristics generally takes longer than learning construction heuristics. This paper focuses 170 on learning construction heuristics for quickly solving instances. 171

172 173

174 175

176

3 MODEL

In this section, we propose a deep learning model with a dynamic encoder and dual-channel decoder structure and provide a detailed introduction to this model.

- 177 178
- 179 3.1 DYNAMIC ENCODER

In a VRP instance with n city nodes, the node feature  $x^i$  for node i, includes 2-dimensional co-181 ordinates and instance-specific features. For example, in TSP, the node coordinates are the node 182 features. The model constructs a complete solution by sequentially selecting nodes. Using the same 183 embeddings at each construction step may lead to poor performance. To address this issue, we use 184 the starting point of the complete solution and the node selected in the previous step to bridge the 185 partial solution with the subproblem composed of all available nodes, complementing the global information. These elements are input into the dynamic encoder along with the available nodes for 187 re-embedding calculations at each step. Since the node from the previous step is input as the starting 188 point of the subproblem, node selection in the current step also serves as a search direction problem 189 for the optimal solution of the previous node. Additionally, incorporating the starting point and the 190 previous node allows the model to dynamically learn the relationship between the partial solution 191 and all available nodes at each step. As the scale of the subproblem input changes at each step, the model learns the relationships between nodes, making it less sensitive to the problem instance scale 192 and thus achieving better generalization. 193

Recalculating node embeddings at each step increases computational cost. Most
 Transformer-based models have several times more encoder attention layers than decoder

196 attention layers. With this structure, at each step, node embeddings must pass through a linear projection layer 197 and multiple attention layers, leading to a significantly high computational cost. To address this, we propose a 199 model with a DEDD structure, as shown in Figure 1. This 200 model includes a dynamic encoder with one attention 201 layer and a dual-channel decoder with multiple attention 202 layers. This structure substantially reduces the computa-203 tional overhead incurred by re-encoding node features in 204 every construction step.

205 To express the node sequence  $X_t$  at step t (where  $t \in$ 206  $\{1, 2, \ldots, n\}$  represents the current construction step), 207 we denote the features of the node selected at step 208 j as  $x_j$ . Thus,  $X_t = (x_1, x_{t-1}, x_t, x_{t+1}, \dots, x_n)$ . 209 The dynamic encoder includes one linear projection 210 layer and one attention layer. For an instance with n211 nodes, the linear projection layer first transforms the 212 node feature sequence  $X_t$  into initial node embeddings  $E^0 = (e_1^0, e_{t-1}^0, e_t^0, e_{t+1}^0, \dots, e_n^0)$ , where  $e_j^0$  denotes the 213 d-dimensional initial embedding of the node selected at 214 step j. These initial embeddings are then processed 215 through an attention layer to obtain the node embeddings



Figure 2: Multi-head Attention Layer

216  $E^1 = (e_1^1, e_{t-1}^1, e_t^1, e_{t+1}^1, \dots, e_n^1)$ . The attention layer consists of a multi-head attention (MHA) sublayer and a feedforward sublayer, as shown in Figure 2. 217 218

Let  $E^{l-1} = (e_1^{l-1}, e_{t-1}^{l-1}, e_t^{l-1}, e_{t+1}^{l-1}, \dots, e_n^{l-1})$  be the input of the *l*-th multi-head attention layer, so the MHA can be defined as follow: 219 220

 $E^{l,m} = \operatorname{Attention}(Q^m, K^m, V^m)$ 

$$Q_j^m, K_j^m, V_j^m = W_Q^m e_j^{l-1}, W_K^m e_j^{l-1}, W_V^m e_j^{l-1}$$
(1)

230

232 233 234

235

236

237

238 239

240 241

242

259

264

221

$$= \operatorname{softmax}(Q^m K^m T / \sqrt{d_k}) V^m, m = 1, 2, \dots, M$$
<sup>(2)</sup>

225 Where  $E^{l,m}$  represents the node embeddings computed by the *m*-th attention head in the *l*-th multi-226 head attention layer. Q, K, V are Query, Key, Value vectors, and  $W_Q^m, W_K^m, W_V^m$  are linear transfor-227 mation matrices of Q, K and V, respectively. The M attention heads respectively perform equations 228 (1) and (2), with  $d_k = d/M$ .

229 After concatenating all  $E^{l,m}$ , they are multiplied by  $W_O$  to obtain the output of the multi-head attention layer. 231

$$Multihead(Q, K, V) = Concat(E^{l,1}, E^{l,2}, \dots, E^{l,M})W_O$$
(3)

$$\hat{e}_j^l = e_j^{l-1} + \text{Multihead}_j(Q, K, V) \tag{4}$$

$$e_j^l = \hat{e}_j^l + FF(\hat{e}_j^l) \tag{5}$$

Finally, using the skip connection layer (He et al., 2016) from equation (4) and the feedforward sublayer composed of two linear projection layers from equation (5), we obtain the output  $e_i^l$  of the self-attention block. The process from equations (1) to (5) is represented as follows:

$$E^{l} = \mathrm{MHA}(E^{l-1}) \tag{6}$$

3.2 DUAL-CHANNEL DECODER

243 To enhance the effectiveness of the node selection strategy at each construction step, we employ 244 a dual-channel output decoder. The first L attention layers of the dual-channel decoder are shared between both channels, while the final attention layer and linear projection layer in each channel 245 have identical structures but do not share parameters. This design choice balances performance 246 improvements with the increased computational cost. Utilizing one attention layer without shared 247 parameters significantly enhances the model's performance with minimal additional computation. 248

249 As the final attention layer and linear projection layer in both channels do not share parameters, the 250 same input may yield different outputs. The randomness from unshared parameters is the primary source of node diversity during the early stages of training, before the model is fully trained. Af-251 ter training, the diversity in node selection strategies primarily stems from the training data of the model. The trained model selects two optimal candidate nodes based on learned strategies at each 253 construction step, rather than randomly choosing nodes of unknown quality, thus more efficiently 254 finding higher-quality solutions. At step t, the encoder's output  $E^1 = (e_1^1, e_{t-1}^1, e_t^1, e_{t+1}^1, \dots, e_n^1)$  is 255 fed into the decoder. The decoder then recalculates the embeddings for  $e_1^1$  and  $e_{t-1}^1$  using two linear 256 projection layers, and concatenates them with  $(e_t^1, e_{t+1}^1, \dots, e_n^1)$  to obtain the input  $\tilde{E}^0$  for its first 257 attention layer. This process is expressed as follows: 258

$$\tilde{E}^0 = \text{Concat}(W_1 e_1^1, W_2 e_{t-1}^1, e_t^1, e_{t+1}^1, \dots, e_n^1)$$
(7)

260 After passing through L attention layers, we obtain the node embeddings  $\tilde{E}^L$ . Then  $\tilde{E}^L$  passes 261 through the last attention layer separately for each of the two channels: 262

$$\tilde{E}^1 = \mathrm{MHA}(\tilde{E}^0) \tag{8}$$

(9)

265 266  $\tilde{E}^L = \mathrm{MHA}(\tilde{E}^{L-1})$ 

267 
$$E^{-} = MHA(E^{-})$$
 (9)  
268  $\tilde{E}_{1}^{L+1} = MHA1(\tilde{E}^{L})$  (10)

$$\tilde{E}_{2}^{L+1} = \mathrm{MHA2}(\tilde{E}^{L}) \tag{11}$$

 $\tilde{E}_r^{L+1}$  is transformed by a linear projection layer into a vector **u**, and finally undergoes a softmax transformation to obtain the selection probability  $\mathbf{p}^t$ , represented as follows:

$$u_i = \begin{cases} W_A \tilde{e}_j^{L+1}, & i \neq 1 \text{ or } 2\\ -\infty, & \text{otherwise} \end{cases}$$
(12)

$$\mathbf{p}^t = \operatorname{softmax}(\mathbf{u}) \tag{13}$$

The high memory and computational costs of the heavyweight decoder structure make it difficult to
use reinforcement learning for training the DEDD model in this paper. Therefore, we use supervised
learning to train the DEDD model. We define the loss function of the DEDD model as follows:

$$loss = -[log(p_1^t) + log(p_2^t)]$$
(14)

Where  $p_r^t$  represents the probability predicted by channel r of the DEDD model for the labeled node in step t.

### 3.3 CANDIDATE NODE SELECTION STRATEGY

The node selection process of the decoder is depicted in Figure 3. After obtaining the selection probabilities  $\mathbf{p}^t$  from both channels, the node with the highest probability in  $\mathbf{p}^t$  is selected as the candidate node for each channel (for instance, in Figure 3, the candidate nodes for instance 1 at step t are 3 and 0). The merits of the two candidate nodes are compared, and the best one is chosen as the final node for step t. During the training process, data is fed into the DEDD model in batches. At each construction step, we calculate the loss for both channels across the entire batch, as shown in Equations (15) and (16), and select the nodes output by the channel with the lower loss as the final nodes for the entire batch. During inference, we calculate the additional distance for both channels across the entire batch, as indicated in Equations (17) and (18), and select the nodes output by the channel with the shorter distance as the final nodes for the entire batch. This approach evaluates the performance of both channels across multiple instances, thereby reducing the randomness associated with single-instance evaluation. 

$$loss_1 = ls_{11} + ls_{12} + \dots + ls_{1b} \tag{15}$$

$$loss_2 = ls_{21} + ls_{22} + \dots + ls_{2b} \tag{16}$$

$$d_1 = d_{11} + d_{12} + \dots + d_{1b} \tag{17}$$

$$d_2 = d_{21} + d_{22} + \dots + d_{2b} \tag{18}$$



Figure 3: The node selection process for the *t*-th construction step.

### 3.4 DESTRUCTION-RECONSTRUCTION APPROACH BASED ON NODE SELECTION STRATEGY

To improve the quality of solutions constructed by the model at once, we adopt a destructionreconstruction (DR) approach based on the node selection strategy. We perform operations on batches of instances. First, a segment, referred to as a partial solution, is randomly extracted from the complete solution and input into the model for step-by-step reconstruction. In each reconstruction step, after the two channels output candidate nodes, the total distance values of all instances in the batch are calculated for both channels. The nodes output by the channel with the smaller total distance is selected as the final nodes for the current step. Once the partial solution is completely reconstructed, the distance of the partial solution before and after reconstruction is compared. If the reconstructed partial solution has a shorter distance, it replaces the original partial solution.

330 331

332

# 4 EXPERIMENT

In this section, we apply the proposed DEDD model to solve TSP and CVRP instances of various sizes. We compare its performance with other learning-based approaches and specialized solvers to assess its effectiveness. The TSP involves finding the shortest path for a salesman who starts from a city, visits all cities exactly once, and returns to the starting city, considering the distances between each pair of cities. The CVRP is a variant of the TSP where vehicles, starting from a depot and each with a capacity limit, aim to satisfy delivery demands for each city.

Our dataset is generated in accordance with the standard data generation procedure established in prior work AM (Kool et al., 2019). We employed the Concorde solver to obtain the optimal solutions for the TSP training set and utilized the HGS to acquire the optimal solutions for the CVRP training set. The training set includes one million TSP and CVRP instances, with sizes ranging from 4 to 100, respectively. The test set comprises 10,000 TSP and 10,000 CVRP instances of size 100, and 128 TSP and 128 CVRP instances of sizes 200, 500, and 1000.

345 Hyperparameter settings: We set the embedding dimension to 128, the encoder to 1 attention layer, 346 the decoder to 5 parameter-shared attention layers, 1 non-shared attention layer. Each attention layer 347 employs 8 heads for multi-head attention, and the feed forward layer dimension is set to 512. The 348 DEDD model is trained separately on one million instances of size 100 for both the TSP and CVRP datasets. For training, we use a batch size of 1024 and the Adam optimizer with a learning rate 349 of  $10^{-4}$ . For the TSP dataset, the learning rate decays by 0.97 per epoch, with training continuing 350 for 150 epochs. For the CVRP dataset, the learning rate decays by 0.9 per epoch, with training 351 lasting for 40 epochs. The DEDD model was trained and tested on an NVIDIA 3090 GPU, with 352 the training to solve TSP instances requiring approximately 8 days and the training to solve CVRP 353 instances requiring approximately 2 days. 354

355 We compare our approach to existing learning-based approaches and classical solvers. Classical solvers include Concorde (Applegate et al., 2006), LKH3 (Helsgaun, 2017), HGS (Vidal, 2022), 356 and OR-Tools (Perron & Furnon, 2019). Learning-based approaches include LEHD (Luo et al., 357 2023), POMO (Kwon et al., 2020), MDAM (Xin et al., 2021), EAS (Hottung et al., 2021b), SGBS 358 (Choo et al., 2022), BQ (Drakulic et al., 2023), and Att-GCN+MCTS (Fu et al., 2021). For optimal-359 ity gap comparison, we use Concorde (Applegate et al., 2006) as the baseline for TSP and LKH3 360 (Helsgaun, 2017) as the baseline for CVRP. Our experiments focus on optimality gaps rather than 361 inference time due to classical solvers running on CPU and learning-based approaches running on 362 GPU, as well as potential differences in programming languages and platforms used for execution.

363 364

# 4.1 EXPERIMENTAL RESULTS

366 Table 1 presents experimental results comparing the DEDD model to various approaches. In TSP, 367 the solutions directly generated by the DEDD model (greedy in Table 1) are of high quality and 368 require minimal inference time. After 100 rounds of DR iteration, the DEDD model's performance surpasses all compared learning-based approaches in Table 1, except for LEHD (Luo et al., 2023). 369 For TSP100, solutions with only 0.009% gap from the baseline are achieved after 100 rounds of DR 370 iteration, nearly matching the baseline. For larger instances like TSP200, TSP500, and TSP1000, 371 which are not in the training set, the DEDD model achieves gaps of approximately 0.03%, 0.2%, and 372 0.8% respectively, after just 300 rounds of DR iteration, closely approaching the baseline results. 373 After 1000 rounds of DR iteration, the DEDD model nearly achieves optimal performance, with gaps 374 of 0.001% for TSP100 and 0.015% for TSP200. The gaps for TSP500 and TSP1000 are 0.142% 375 and 0.611%, respectively. 376

In CVRP, the DEDD model demonstrates superior performance relative to TSP. Due to the memory constraints of a single NVIDIA 3090 GPU, 10,000 CVRP100 instances are divided into two batches

| 515        | Table 1. The experin                  | nemai res      | unts 101 |                | VICI UI | uuci unnoi     | inity uist | induce ms     | tances. |
|------------|---------------------------------------|----------------|----------|----------------|---------|----------------|------------|---------------|---------|
| 380        | Method                                | TSP100         |          | TSP200         |         | TSP500         |            | TSP1000       |         |
| 381<br>382 | i i i i i i i i i i i i i i i i i i i | Gap            | Time     | Gap            | Time    | Gap            | Time       | Gap           | Time    |
| 383        | Concorde                              | 0.000%         | 34m      | 0.000%         | 3m      | 0.000%         | 32m        | 0.000%        | 7.8h    |
| 384        | LKH3                                  | 0.000%         | 56m      | 0.000%         | 4m      | 0.000%         | 32m        | 0.000%        | 8.2h    |
| 385        | OR-Tools                              | 2.368%         | 11h      | 3.618%         | 17m     | 4.682%         | 50m        | 4.885%        | 10h     |
| 396        | Att-GCN+MCTS                          | 0.037%         | 15m      | 0.884%         | 2m      | 2.536%         | 6m         | 3.223%        | 13m     |
| 300 -      | MDAM bs50                             | 0.388%         | 21m      | 1.996%         | 3m      | 10.065%        | 11m        | 20.375%       | 44m     |
| 307        | POMO augx8                            | 0.134%         | 1m       | 1.533%         | 5s      | 22.187%        | 1m         | 40.570%       | 8m      |
| 388        | SGBS                                  | 0.060%         | 40m      | 0.562%         | 4m      | 11.550%        | 54m        | 26.035%       | 7.4h    |
| 389        | EAS                                   | 0.057%         | 6h       | 0.496%         | 28m     | 17.080%        | 7.8h       | -             | -       |
| 390        | BQ greedy                             | 0.579%         | 0.6m     | 0.895%         | 3s      | 1.834%         | 0.4m       | 3.965%        | 2.4m    |
| 391        | BQ bs16                               | 0.046%         | 11m      | 0.224%         | 1m      | 0.896%         | 6m         | 2.605%        | 38m     |
| 392        | LEHD RRC 1000                         | 0.002%         | 2.1h     | 0.020%         | 12.3m   | 0.182%         | 1.3h       | 0.745%        | 7.2h    |
| 393        | DEDD greedy                           | 0.533%         | 0.6m     | 0.870%         | 0.07m   | 1.714%         | 0.4m       | 2.899%        | 2.0m    |
| 394        | DEDD DR 50                            | 0.182%         | 10.6m    | 0.118%         | 0.8m    | 0.454%         | 6.9m       | 1.245%        | 31.5m   |
| 395        | DR 100                                | 0.009%         | 19.0m    | 0.061%         | 1.8m    | 0.325%         | 12.5m      | 1.069%        | 1.0h    |
| 396        | DR 300                                | 0.004%         | 51.4m    | 0.028%         | 5.1m    | 0.198%         | 34.3m      | 0.801%        | 2.9h    |
| 397        | DR 500                                | 0.002%         | 1.4h     | <u>0.019%</u>  | 8.4m    | <u>0.170%</u>  | 53.9m      | <u>0.688%</u> | 4.7h    |
| 308        | DR 1000                               | 0.001%         | 2.8h     | 0.015%         | 17.0m   | 0.142%         | 1.8h       | 0.611%        | 9.4h    |
| 300        |                                       | CVRP100        |          | CVRP200        |         | CVRP500        |            | CVRP1000      |         |
|            | LKH3                                  | 0.000%         | 12h      | 0.000%         | 2.1h    | 0.000%         | 5.5h       | 0.000%        | 7.1h    |
| 400        | HGS                                   | -0.533%        | 4.5h     | -1.126%        | 1.4h    | -1.794%        | 4h         | -2.162%       | 5.3h    |
| 401        | OR-Tools                              | 6.193%         | 2h       | 6.894%         | 1h      | 9.112%         | 2.2h       | 11.662%       | 3h      |
| 402        | MDAM bs50                             | 2.211%         | 25m      | 4.304%         | 3m      | 10.498%        | 12m        | 27.814%       | 47m     |
| 403        | POMO augx8                            | 0.689%         | 1m       | 4.866%         | 7s      | 19.901%        | 1m         | 128.89%       | 10m     |
| 404        | SGBS                                  | 0.079%         | 40m      | 2.581%         | 1m      | 15.343%        | 16m        | 136.98%       | 2.3h    |
| 405        | EAS                                   | -0.234%        | 15h      | 0.640%         | 33m     | 11.042%        | 9.3h       | -             | -       |
| 406        | BQ greedy                             | 2.993%         | 0.7m     | 3.527%         | 4s      | 5.121%         | 0.4m       | 9.812%        | 2.4m    |
| 407        | BQ bs16                               | 0.611%         | 10m      | 1.141%         | 0.6m    | 2.991%         | 6m         | 7.784%        | 39m     |
| 408        | LEHD RRC 1000                         | -0.100%        | 2.7h     | <u>-0.346%</u> | 14.5m   | -0.01%         | 1.35h      | 2.484%        | 7.8h    |
| 409        | DEDD greedy                           | 3.557%         | -        | 3.055%         | 0.1m    | 2.877%         | 0.4m       | 6.456%        | 2.2m    |
| 410        | DEDD DR 50                            | 0.525%         | -        | 0.402%         | 1.4m    | 0.948%         | 6.0m       | 4.149%        | 38.7m   |
| 411        | DR 100                                | 0.257%         | -        | 0.123%         | 2.5m    | 0.606%         | 10.4m      | 3.602%        | 1.1h    |
| /110       | DR 300                                | 0.007%         | -        | -0.196%        | 6.8m    | 0.136%         | 32.8m      | 2.688%        | 3.0h    |
| 410        | DR 500                                | -0.071%        | -        | -0.302%        | 11.5m   | <u>-0.015%</u> | 57.8m      | <u>2.332%</u> | 4.9h    |
| 413        | DR 1000                               | <u>-0.148%</u> | 3.5h     | -0.446%        | 22.6m   | -0.245%        | 1.9h       | 1.879%        | 10.0h   |
| 414 .      |                                       |                |          |                |         |                |            |               |         |

Table 1: The experimental results for TSD and CVDD under uniformly distributed instances

415

417

418

378

070

416 of 5000 each for inference, and the results are averaged. Table 1 presents the total time after completing 1000 rounds of DR iteration for both batches. The solutions directly generated by DEDD show gap of about 3% from the baseline for CVRP100, CVRP200, and CVRP500, and about 6% 419 for CVRP1000.

420 After 300 rounds of DR iteration, solution quality significantly improves. Specifically, the gap 421 for CVRP100 is 0.007%, CVRP200's solution surpasses LKH3 (Helsgaun, 2017) by -0.196%, 422 CVRP500's gap is close to the baseline at 0.136%, and CVRP1000 achieves a gap of 2.688%. 423

After 500 rounds of DR iteration, the DEDD model clearly outperforms LKH3 (Helsgaun, 2017) on 424 CVRP100, CVRP200, and CVRP500. After 1000 rounds of DR iteration, the gap for CVRP1000 is 425 also within 2%. 426

427 As shown in Table 1, the DEDD model's performance slightly fell short of EAS only in the 428 CVRP100 solution. After 300 rounds of DR iterations, the DEDD model surpass five learning-based methods (except EAS and LEHD) and the traditional solver OR-Tools across four scales. For larger 429 scales, such as CVRP200, CVRP500, and CVRP1000, the DEDD model demonstrate stronger ca-430 pabilities, surpassing six learning-based methods (except LEHD) and OR-Tools with only 50 rounds 431 of DR iterations. These results indicate that the DEDD model performs excellently.

### 432 5 **ABLATION STUDY**

433 434

460 461

We assess the effectiveness of various model components using TSP instances across four scales. 435 Specifically, we compare the static encoder-dual channel decoder (SEDD) and the dynamic encoder-436 single channel decoder (DESD) with the dynamic encoder-dual channel decoder (DEDD) to evaluate 437 the individual contributions of the dynamic encoder and dual-channel decoder to the model's per-438 formance. The dynamic encoder embeds nodes for subproblems at each construction step, while 439 the static encoder embeds all nodes in a single initial step. The dual-channel decoder generates 440 two probability matrices, from which higher-quality nodes are selected based on a node selection strategy, while the single-channel decoder outputs only one probability matrix. Table 2 shows that 441 DEDD outperforms DESD and SEDD in all metrics for TSP100. However, for TSP200, TSP500, 442 and TSP1000, the solutions output by DEDD are slightly inferior to those of DESD or SEDD. 443

444 First, the DEDD solutions may be slightly inferior to those of SEDD because models of two struc-445 tures are trained on a dataset of one million TSP instances ranging from size 4 to 100, thus learning relationships among nodes in small-scale instances. When constructing complete solutions for 446 TSP100, dynamic encoders improve performance comprehensively, as models of all three struc-447 tures are trained on problems of related sizes. However, for larger-scale problems, neither dynamic 448 nor static encoders learn to encode features of nodes at this scale during training, complicating the 449 comparison of encoding approaches, as shown in the greedy metrics in Table 2. 450

451 Nonetheless, improvements can be made to the initial solution after it is directly output by the model. We continuously improve solution quality using the DR strategy. Since partial solution 452 problems are present in the training set, the advantages and disadvantages of dynamic and static 453 encoders are apparent. Specifically, a more suitable encoder can greatly improve partial solutions, thereby enhancing the overall solution quality. Table 2 shows that after multiple rounds of DR 455 iteration, models using dynamic encoders (DEDD) outperform those using static encoders (SEDD) 456 on all problem scales, proving the effectiveness of dynamic encoders. 457

Secondly, DEDD solutions may be slightly inferior to those of DESD for the following reasons. The 458 dual-channel decoder offers richer node selection for each construction 459

462 463 **TSP100** 464 DR 50 DR 100 DR 300 DR 500 DR 1000 greedy 465 DESD gap 0.554% 0.023% 0.013% 0.0042% 0.0024% 0.0013% 466 SEDD gap 0.542% 0.033% 0.014% 0.0041% 0.0023% 0.0013% 467 DEDD gap 0.533% 0.018% 0.009% 0.0040% 0.0022% 0.0012% 468 469 **TSP200** 470 DR 100 DR 300 greedy DR 50 DR 500 DR 1000 471 DESD gap 0.922% 0.121% 0.075% 0.033% 0.023% 0.0155% 472 SEDD gap 0.826% 0.119% 0.068% 0.029% 0.021% 0.0152% 473 DEDD gap 0.870% 0.118% 0.061% 0.028% 0.019% 0.0151% 474 **TSP500** 475 476 greedy DR 50 DR 100 DR 300 DR 500 DR 1000 477 DESD gap 1.684% 0.474% 0.317% 0.210% 0.188% 0.165% 478 SEDD gap 1.732% 0.455% 0.329% 0.224% 0.190% 0.163% 479 DEDD gap 1.714% **0.454%** 0.325% 0.198% 0.170% 0.142% 480 TSP1000 481 482 greedy DR 50 DR 100 DR 300 DR 500 DR 1000 483 DESD gap 2.713% 1.155% 1.081% 0.828% 0.693% 0.616% 484 SEDD gap 3.091% 1.381% 1.181% 0.944% 0.726% 0.837% 485 DEDD gap 2.899% 1.245% 1.069% 0.801% 0.688% 0.611%

Table 2: Experimental results of models with different architectures on uniformly distributed instances of TSP.

step, but this does not guarantee a higher-quality complete solution. Thus, the results of a single
solve may exhibit some randomness and yield slightly lower-quality solutions. However, after sufficient rounds of reconstructing partial solutions, randomness is greatly reduced, resulting in highquality solutions. Table 2 shows that after 1000 rounds of DR, the quality of DEDD's solutions
surpasses DESD, demonstrating the effectiveness of the dual-channel decoder structure.

Finally, we comprehensively compare the experimental results of the three structures. Since supervised learning-trained encoders may not accurately encode node features for problem scales outside
the training set, and the solutions produced by dual-channel decoders in one solve may exhibit some randomness, we focus on the performance of models after multiple rounds of DR in the ablation experiment. Table 2 shows that after multiple rounds of DR, DEDD performs best among the three model structures, proving its effectiveness.

497 498

499

508

523

524

525

526

527

528

529

530

531

6 CONCLUSION

500 This paper introduces a novel model with dynamic encoders and a dual-channel decoder to learn 501 construction heuristic for routing problems. The model uses supervised learning and incorporates 502 dynamic encoders and a subproblem decomposition strategy tailored to routing problems, enhanc-503 ing performance. The dual-channel decoder enables rich node selection and customized strategies, 504 further enhancing performance. Experimental results show that the DEDD model performs excep-505 tionally well in both TSP and CVRP, with particularly strong results in CVRP. Future work will focus on improving cooperative decision-making between decoder channels and exploring the potential of 506 training the DEDD model with reinforcement learning. 507

# 509 REFERENCES

- 510 511 David Applegate, Ribert Bixby, Vasek Chvatal, and William Cook. Concorde tsp solver, 2006.
- Irwan Bello, Hieu Pham, Quoc V Le, Mohammad Norouzi, and Samy Bengio. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.
- 515 Dimitri Bertsekas. Dynamic programming and optimal control. belmont, ma: Athena scientific.
   516 Appendix E, 2000.
- Xinyun Chen and Yuandong Tian. Learning to perform local rewriting for combinatorial optimization. *Advances in neural information processing systems*, 32, 2019.
- Jinho Choo, Yeong-Dae Kwon, Jihoon Kim, Jeongwoo Jae, André Hottung, Kevin Tierney, and
   Youngjune Gwon. Simulation-guided beam search for neural combinatorial optimization. Advances in Neural Information Processing Systems, 35:8760–8772, 2022.
  - Paulo R d O Costa, Jason Rhuggenaath, Yingqian Zhang, and Alp Akcay. Learning 2-opt heuristics for the traveling salesman problem via deep reinforcement learning. In *Asian conference on machine learning*, pp. 465–480. PMLR, 2020.
  - Michel Deudon, Pierre Cournut, Alexandre Lacoste, Yossiri Adulyasak, and Louis-Martin Rousseau. Learning heuristics for the tsp by policy gradient. In *Integration of Constraint Programming, Artificial Intelligence, and Operations Research: 15th International Conference, CPAIOR 2018, Delft, The Netherlands, June 26–29, 2018, Proceedings 15*, pp. 170–181. Springer, 2018.
- Darko Drakulic, Sofia Michel, Florian Mai, Arnaud Sors, and Jean-Marc Andreoli. Bq-nco:
   Bisimulation quotienting for generalizable neural combinatorial optimization. *arXiv preprint arXiv:2301.03313*, 2023.
- Matteo Fischetti, Paolo Toth, and Daniele Vigo. A branch-and-bound algorithm for the capacitated vehicle routing problem on directed graphs. *Operations Research*, 42(5):846–859, 1994.
- Zhang-Hua Fu, Kai-Bin Qiu, and Hongyuan Zha. Generalize a small pre-trained model to arbitrarily
   large tsp instances. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 7474–7482, 2021.

| 540<br>541<br>542        | A Hanif Halim and IJAoCMiE Ismail. Combinatorial optimization: comparison of heuristic algorithms in travelling salesman problem. <i>Archives of Computational Methods in Engineering</i> , 26: 367–380, 2019.  |  |  |  |  |  |
|--------------------------|---|--|--|--|--|--|
| 543<br>544<br>545        | Laiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-<br>nition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp.   |  |  |  |  |  |
| 546<br>547               | 770–778, 2016.<br>Keld Helsgaun. An extension of the lin-kernighan-helsgaun tsp solver for constrained traveling  |  |  |  |  |  |
| 548<br>549               | salesman and vehicle routing problems. <i>Roskilde: Roskilde University</i> , 12:966–980, 2017.   |  |  |  |  |  |
| 550<br>551<br>552        | 1997.   |  |  |  |  |  |
| 553<br>554               | John J Hopfield and David W Tank. "neural" computation of decisions in optimization problems.<br><i>Biological cybernetics</i> , 52(3):141–152, 1985.   |  |  |  |  |  |
| 555<br>556<br>557        | André Hottung and Kevin Tierney. Neural large neighborhood search for the capacitated vehicle routing problem. In <i>ECAI 2020</i> , pp. 443–450. IOS Press, 2020.  |  |  |  |  |  |
| 558<br>559<br>560        | André Hottung, Bhanu Bhandari, and Kevin Tierney. Learning a latent search space for routing prob-<br>lems using variational autoencoders. In <i>International Conference on Learning Representations</i> , 2021a.  |  |  |  |  |  |
| 561<br>562<br>563        | André Hottung, Yeong-Dae Kwon, and Kevin Tierney. Efficient active search for combinatorial optimization problems. <i>arXiv preprint arXiv:2106.05126</i> , 2021b.  |  |  |  |  |  |
| 564<br>565               | Chaitanya K Joshi, Thomas Laurent, and Xavier Bresson. An efficient graph convolutional network technique for the travelling salesman problem. <i>arXiv preprint arXiv:1906.01227</i> , 2019.   |  |  |  |  |  |
| 567<br>568<br>569        | Chaitanya K Joshi, Quentin Cappart, Louis-Martin Rousseau, and Thomas Laurent. Learning the travelling salesperson problem requires rethinking generalization. <i>Constraints</i> , 27(1):70–98, 2022.  |  |  |  |  |  |
| 570<br>571               | Elias Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, and Le Song. Learning combinatorial opti-<br>mization algorithms over graphs. <i>Advances in neural information processing systems</i> , 30, 2017.  |  |  |  |  |  |
| 573<br>574<br>575        | Wouter Kool, Herke van Hoof, Joaquim Gromicho, and Max Welling. Deep policy dynamic pro-<br>gramming for vehicle routing problems. In <i>International conference on integration of constraint</i><br><i>programming, artificial intelligence, and operations research</i> , pp. 190–213. Springer, 2022. |  |  |  |  |  |
| 576<br>577<br>578        | Yeong-Dae Kwon, Jinho Choo, Byoungjip Kim, Iljoo Yoon, Youngjune Gwon, and Seungjai Min.<br>Pomo: Policy optimization with multiple optima for reinforcement learning. <i>Advances in Neural</i><br><i>Information Processing Systems</i> , 33:21188–21198, 2020.   |  |  |  |  |  |
| 580<br>581               | Eugene L Lawler and David E Wood. Branch-and-bound methods: A survey. <i>Operations research</i> , 14(4):699–719, 1966.   |  |  |  |  |  |
| 582<br>583<br>584<br>585 | Bingjie Li, Guohua Wu, Yongming He, Mingfeng Fan, and Witold Pedrycz. An overview and experimental study of learning-based optimization algorithms for the vehicle routing problem. <i>IEEE/CAA Journal of Automatica Sinica</i> , 9(7):1115–1138, 2022.  |  |  |  |  |  |
| 586<br>587               | Hao Lu, Xingwen Zhang, and Shuang Yang. A learning-based iterative method for solving vehicle routing problems. In <i>International conference on learning representations</i> , 2019.  |  |  |  |  |  |
| 588<br>589<br>590<br>591 | Fu Luo, Xi Lin, Fei Liu, Qingfu Zhang, and Zhenkun Wang. Neural combinatorial optimization with heavy decoder: Toward large scale generalization. <i>Advances in Neural Information Processing Systems</i> , 36:8845–8864, 2023.  |  |  |  |  |  |
| 592<br>593               | Qiang Ma, Suwen Ge, Danyang He, Darshan Thaker, and Iddo Drori. Combinatorial opti-<br>mization by graph pointer networks and hierarchical reinforcement learning. <i>arXiv preprint</i><br><i>arXiv:1911.04936</i> , 2019.   |  |  |  |  |  |

- 594 Yining Ma, Jingwen Li, Zhiguang Cao, Wen Song, Le Zhang, Zhenghua Chen, and Jing Tang. Learning to iteratively solve routing problems with dual-aspect collaborative transformer. Ad-596 vances in Neural Information Processing Systems, 34:11096–11107, 2021. 597 Mohammadreza Nazari, Afshin Oroojlooy, Lawrence Snyder, and Martin Takác. Reinforcement 598 learning for solving the vehicle routing problem. Advances in neural information processing systems, 31, 2018. 600 601 Vincent Laurent Perron and Furnon. Google's or-tools. 602 https://developers.google.com/optimization, 2019. Accessed: 2024-02-603 16. 604 Moshe Sniedovich. Dynamic programming: foundations and principles. CRC press, 2010. 605 606 Paolo Toth and Daniele Vigo. Vehicle routing: problems, methods, and applications. SIAM, 2014. 607 Peter JM Van Laarhoven, Emile HL Aarts, and Jan Karel Lenstra. Job shop scheduling by simulated 608 annealing. Operations research, 40(1):113–125, 1992. 609 610 Vijay V Vazirani. Approximation algorithms, 2001. 611 Matthew Veres and Medhat Moussa. Deep learning for intelligent transportation systems: A survey 612 of emerging trends. IEEE Transactions on Intelligent transportation systems, 21(8):3152–3168, 613 2019. 614 615 Thibaut Vidal. Hybrid genetic search for the cvrp: Open-source implementation and swap\* neigh-616 borhood. Computers & Operations Research, 140:105643, 2022. 617 Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. Advances in neural informa-618 tion processing systems, 28, 2015. 619 620 Christopher JCH Watkins and Peter Dayan. Q-learning. Machine learning, 8:279–292, 1992. 621 J Wesley Barnes and MANUEL LAGUNA. Solving the multiple-machine weighted flow time prob-622 lem using tabu search. *IIE transactions*, 25(2):121–128, 1993. 623 624 David P Williamson and David B Shmoys. The design of approximation algorithms. Cambridge 625 university press, 2011. 626 Yaoxin Wu, Wen Song, Zhiguang Cao, Jie Zhang, and Andrew Lim. Learning improvement heuris-627 tics for solving routing problems. IEEE transactions on neural networks and learning systems, 628 33(9):5057-5069, 2021. 629 630 Liang Xin, Wen Song, Zhiguang Cao, and Jie Zhang. Multi-decoder attention model with embedding 631 glimpse for solving vehicle routing problems. In Proceedings of the AAAI Conference on Artificial 632 Intelligence, volume 35, pp. 12042–12049, 2021. 633 Emre Yolcu and Barnabás Póczos. Learning local search heuristics for boolean satisfiability. Ad-634 vances in Neural Information Processing Systems, 32, 2019. 635 636 637 638 639 640 641 642 643 644 645
- 646
- 647