# Learning Enhanced Optimisation for Routing Problems

**Nasrin Sultana, Jeffrey Chan, Tabinda Sarwar, Babak Abbasi, A. K. Qin**

September 20, 2021

## Abstract

Deep learning approaches have shown promising results in solving routing problems. However, there is still a substantial gap in solution quality between machine learning and operations research algorithms. Recently, another line of research has been introduced that fuses the strengths of machine learning and operational research algorithms. In particular, search perturbation operators have been used to improve the solution. Nevertheless, using the perturbation may not guarantee a quality solution. This paper presents "Learning to Guide Local Search" (L2GLS), a learning-based approach for routing problems that uses a penalty term and reinforcement learning to adaptively adjust search efforts. L2GLS combines local search (LS) operators' strengths with penalty terms to escape local optimals. Routing problems have many practical applications, often presetting larger instances that are still challenging for many existing algorithms introduced in the learning to optimise field. We show that L2GLS achieves the new state-of-the-art results on larger TSP and CVRP over other machine learning methods.

***K*eywords** Capacitated Vehicle Routing · Travelling Salesman Problem · Deep Learning · Learning to Optimise

## 1 Introduction

Routing problems are an important class of combinatorial optimisation problems (COPs), which have many real-life applications, e.g., supply chain and warehouse management, aviation planning, healthcare scheduling, and hardware design [1]. Among the routing problems, the Travelling Salesman Problem (TSP) and Vehicle Routing Problem (VRP) are the two most prevalent ones. TSP is defined as finding a tour of all cities/locations where each city is visited only once and has a minimum total distance travelled. The VRP is to find the optimal set of routes for a fleet of vehicles. VRP has many variants; one example is Capacitated vehicle routing problems (CVRP) that aims to find a set of routes with minimum cost to fulfil customers' demands without violating the vehicle capacity constraints.

Routing problems [2] can be solved efficiently with many optimisation methods, including Branch-and-Bound [3], Local Search (LS) [4], Lagrangian Relaxation [5] and Tabu Search [6]. A highly specialised algorithm was designed for TSP, namely, Concorde [7], which is widely used as an exact TSP solver for large instances. Nevertheless, solving larger instances is difficult for Concorde because of the exponential growth in execution time with increasing instance size. Another approach is manually designed heuristics that can discover close-to-optimal results [8] [9]. For example for CVRP, LKH3 [8] is widely used and is a penalty-function-based extension of the classical Lin-Kernighan heuristic [8]. These two-classic algorithms [7][8] are able to provide a solution with an average cost of 7.77 for 100 cities TSP and 15.65 for 100 customers CVRP but with a long-running time. For instance, for 100 nodes TSP instances, Concorde [10] needs on average one hour to solve one instance, and for 100 nodes CVRP instances, LKH3 needs on average 13 hours to solve an instance.

As highlighted, the traditional solvers have many issues. Recently, Deep Learning (DL) is used to learn heuristics for solving routing problems [11] [12] [13] automatically. They typically develop DL and Reinforcement Learning (RL) frameworks to predict a complete solution from scratch. Their frameworks have shown their competency to obtain an approximate solution to solve RPs. Once the learning-based methods are trained, the execution time is fast. However, current learning-based methods have shortcomings. First, these learning-based methods are designed to act as a constructive heuristic, such as constructing a solution to the TSP sequence one node at a time. However, this model might not be scalable when there are a larger number of nodes (customers). Second, learning-based methods rely on gradient information to guide the search (greedy search or beam search), which may not be available because the solution space is undifferentiable or finding a differentiable surrogate is difficult. In addition, another issue is that most of the work in the Learning to Optimise (L2O) field evaluated their models on randomly generated and easy to solve TSP and CVRP instances and focused on solving instances up to 100 nodes. Scaling to large and real-world instances is still an open question. Solving larger problems is highly time-consuming for most of the current operation research algorithms [8] [14], and [15]. Therefore, all approaches including exact, heuristics and learning-based, have strengths and weaknesses. Naturally, one may want to combine the strengths of operation research and learning-based algorithms to build a mechanism for solving large routing problems efficiently.

There have been some previous approaches based on operation research algorithms and learning-based heuristics and tested on TSP [16] and CVRP [17] [18] [19]. Hottung et al. [19] used destroy operators to improve solutions, while L2I [17] used perturbation to escape local minima and combined with search operators [17] to improve solutions. However, perturbation operators have some issues, such as they do not guarantee an improvement in the quality of the returned solution [20] for all the problems. L2I [17] demonstrated that when the magnitude of perturbation is too large, the resulting solution generally becomes much worse, and the algorithm will take many improvement steps to repair the deterioration.

To address this challenge, we design a method that can scale to large instances that are quite common in real-world problems. In our work, instead of using perturbation operators to avoid local minimum, our method augments the cost function to include a set of penalty terms [21] and passes the new modification to LS. The goal is to escape the local minimum and guide the LS process to distribute search efforts. Furthermore, previous works never investigated which combination of search operators was successful for various routing problems in the L2O field. The current literature focuses on algorithm optimisation performance rather than a thorough examination of various operators. Since selected operators have an impact on the algorithm's efficiency, we also need to analyse them thoroughly. Hence, we analysed many combinations of LS operators and constructed a rich set of LS operators to improve the solution [22] [23] [24]. One of the important observations is that successful heuristics for the route optimisation problems do not necessarily require new operators [25]. We investigate multiple operators and their combinations for the case of LS here. We consider a group of operators, namely 2-opt [4], relocate [24], three permutations [22] and location swap [26] (all operators are described in Section 3.4). The idea behind our method is to find a better solution through search operations by taking improvement progressively into account, learning using RL instead of handcrafted heuristics.

Our proposed approach can improve the initial solution, progressively guiding the search process using penalty terms among feasible solutions using LS operator, followed by a reinforcement-based manager to learn the set of LS operators that are more effective for the RP. Our main contributions are as follows:

- We propose a scalable learning-based algorithm that solve large scale routing problems, achieving new state-of-the-art results. Proposed model can generalise on the benchmark datasets, such as TSPLIB [27] and for CVRP, mentioned in Uchoa et al. [28]. Our algorithm outperforms Concorde on TSP in terms of computation time and solution quality for randomly generated instances, that has gained attention in the last few years for L2O algorithms.

- We propose a framework, where LS operators place in action. LS uses a set of penalty terms when stuck in local minima and searches the promising areas of the search space, i.e., instead of changing the search direction, our method dynamically changing the objective.

- The recent line of research using LS operators shows the potential to solve routing problems using RL. Nevertheless, they used search operators without investigating which combination of operators are effective. We studied a set of LS operators and a learning approach that automatically finds an effective set of operators and the order to apply for routing problems.

## 2   Related Work

DL and RL have recently been proposed to solve COPs. Learning-based method for routing problems can be categorised as constructive [29] [12] [13] [30] [11] [31] [32] [33] [34] and improvement heuristics [18] [19] [35] [17] [16] [36]. Constructive methods tend to predict the next node given a partially constructed tour and build up to a solution node by node. Improvement heuristics improve a solution by iteratively performing search based on LS heuristics to improve the solution [37], where a neural network guides a LS algorithm that iteratively finds promising solutions.

A pioneering work in this area is Pointer Networks (PN) [29] based on attention mechanism and a supervised learning approach to solve TSP. It relies on having an existing set of good solutions or a good solver to generate training instances. One recent paper used a supervised learning approach to train a small-scale model, which could be repetitively used to build heat maps for TSP instances of arbitrarily large instances based on graph sampling, graph converting and heat maps merging techniques. Then the heat maps are fed into a reinforcement learning approach to guide the search for high-quality solutions based on the Monte Carlo tree search [38]. In [13], a constructive neural method was proposed that uses a recurrent neural network (RNN) decoder and an attention mechanism to build solutions for the CVRP and uses an actor-critic RL for training. A graph attention network is used in the method called AM [11] and generates solutions for different routing problems trained via RL, including TSP and CVRP. They trained their model using policy gradient RL with a baseline based on a deterministic greedy roll-out. ERRL [32] is another RL based approach that introduces more exploration via stochastic policies. POMO [33] introduced an end-to-end approach for building a heuristic solver based on policy optimisation with multiple optima. Bresson et al. [34] proposed adapting the popular Transformer architecture to solve TSP. Training is carried out through RL (without TSP training solutions) and decoding uses beam search. Other lines of research are Deep Policy Dynamic Programming that aims to combine the strengths of learned neural heuristics with those of dynamic programming algorithms [39] and the improvement methods [18] [19] [35] and [17]. Many improvement methods have been developed recently [19] [18] [16] [35] [17]. LHI [16] learns which combinations of local operators to apply to solve RPs. RL2OPT [35] proposed a deep RL approach to learn an LS heuristic based on 2-opt operators. Chen et al. [18] focuses on improving heuristics and proposes an RL based improvement approach that chooses a region of a graph representation of the problem and then selects and applies established local heuristics. L2I [17] proposed a combination of LS as an improvement operator and ML algorithm. To escape the local minimum, they further improved the solution by a perturbation operator in L2I [17].

LS methods can easily get stuck in local minimum [40]. In the OR field, various meta-heuristics have been proposed to avoid local optima. They show promising results to solve RPs [41] [42]. However, meta-heuristics still require expert knowledge. Our work is similar to L2I [17], but differs in two important aspects; we use a different set of improving LS operators and we do not use perturbation operators to avoid local minimum. Hence, to produce promising results, we refine LS operators for routing problems and to tackle the local minima problems, we propose to include a set of penalty terms [21] that will adaptively guide the LS operators.

## 3   Preliminaries

In this section we first describe notation used in our approach, and we provide a formal definition of TSP, CVRP, and LS operators used in our approach.

### 3.1   Notations

### 3.2   TSP

Let there be $N$ nodes, each representing a city, and a matrix $D = [d_{i,j}]$, which gives the distance between two cities $i$ and $j$. The goal in TSP is to find a sequence of cities (tour) that visits each city exactly once and the total inter-city distance is of minimum length. A tour can be represented as a permutation $\pi$ on the $N$ cities if we interpret $\pi(i)$ to be the next city visited after city $i$ in the tour, $i = 1, \cdots\cdots, N$, and $\pi(N) = 1$. The objective of TSP can be written as:

$$L(\pi) = \sum_{i=1}^{N} d_{i,\pi(i)} \tag{1}$$

Table 1: Summary of Symbols.

| Symbols | Definition |
| --- | --- |
| n | Nodes |
| m | Coordinates of nodes |
| $\pi$ | Permutation/Solution |
| $L(\pi)$ | Objective function |
| $\pi_L$ | Local minimum |
| $h(\pi)$ | Augmented Objective function |
| $M$ | Maximum roll out steps until |
| $I$ | Improvement steps |
| S | State |
| A | Action |
| G_i | Demand for customer sites for CVRP |
| C_i | Capacity of vehicles for CVRP |
| $p_\theta(A|S)$ | Probability distribution |

### 3.3 CVRP

In CVRP, there is a depot and a set of N customers. Each customer $i, i \in 1; \cdots; N$, has a demand $G_i$ to be satisfied. There is a vehicle that makes a number of trips from the depot and serves a number of customers until the vehicle's total demand exceeds the vehicle's capacity $C$, from which the vehicle has to return to the depot. Similar to the TSP problem, the distance between two customers $i$ and $j$ is denoted by $[d_{ij}]$. The traveling cost $[d_{ij}]$ is the cost of a vehicle going from node i to j, with $i; j \in V = 0; 1; \cdots; N$. Here, the depot is denoted by node 0 for convenience. The goal is to schedule a number of trips that has the minimal total trip distance to serve all customers and respecting the vehicle capacity constraints for each trip and can be mathematically written as:

$$min \sum_{i \in V} \sum_{j \in V} d_{ij} x_{ij}$$

subject to,

$$\sum_{i \in V} x_{i,j} = 1 \ \forall_j \in V \setminus \{0\} \cdots\cdots\cdots\cdots 1,$$

$$\sum_{j \in V} x_{i,j} = 1 \ \forall_i \in V \setminus \{0\} \cdots\cdots\cdots\cdots 2,$$

$$\sum_{i \in V} x_{i,0} = K \cdots\cdots\cdots\cdots 3,$$

$$\sum_{j \in V} x_{0,j} = K \cdots\cdots\cdots\cdots 4,$$

$u_i - u_i + Cx_{ij} \le C - G_j, \forall_{i,j} \in V \ \{0\}, i \ne j$
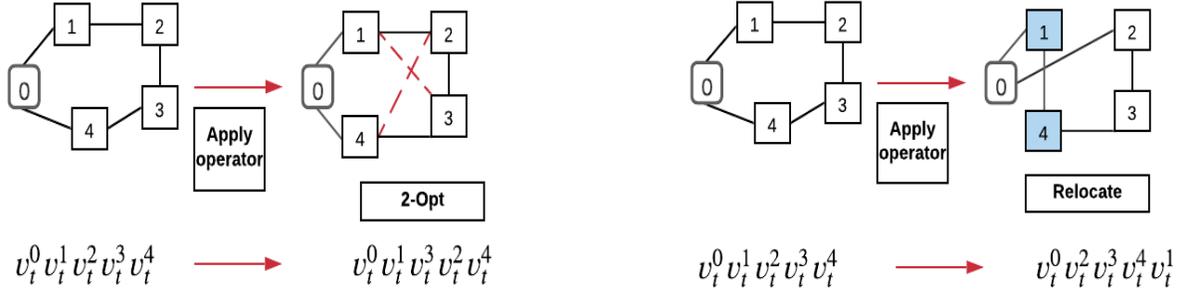s.t $G_i + G_j \le C \cdots\cdots\cdots\cdots 5$

$G_i \le u_i \le C, \forall_i \in V \setminus \{0\} \cdots\cdots\cdots\cdots 6$

$x_{ij} \in 0, 1, \forall_{ij} \in V$

Where K is the number of vehicles available (without loss of generality, it is assumed that K = N for the CVRP we consider), constraints (1) and (2) specify that each customer is visited exactly once, while constraints (3) and (4) specify the in and out degree of the depot, respectively. Constraints (5) and (6) implies the vehicle capacity requirements.
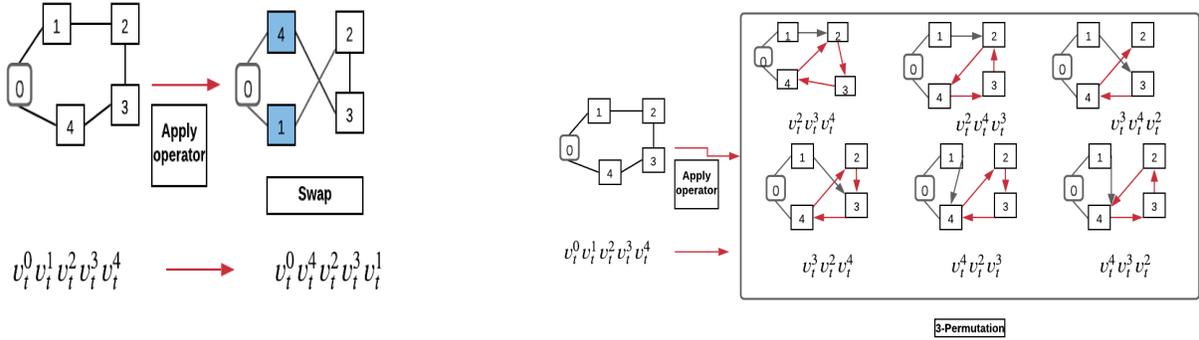
### 3.4 Local Search Operators

LS improve feasible solutions through a search procedure, i.e., it starts with an initial feasible solution and replaces a previous solution with a more optimal solution. Figure 1 gives an illustration of LS operators used in our model. The details of operators applied on routing problems are as follows:

(a) After applying operator 2-Opt, route changes to a new solution with red dashed lines.

(b) After applying relocate operator, location puts one location after another location $(v_1, v_4)$.

(c) After applying swap operator, location exchanges $(v_1, v_4)$

(d) Operator: Three Permutation

Figure 1: An illustration of LS operators.

### 3.4.1 2-opt

Croes et el [4] first introduced the 2-optimisation (2-opt) method, which is a simple and commonly used operator. The idea of 2-opt is to exchange the links between two pairs of subsequent nodes. Figure 1a depicts an example of 2-opt.

### 3.4.2 Relocate

[24] is the process where a selected node (target) is moved from its current position in the tour to another position (destination). Hence, the position of the selected node is relocated. Each relocation of a node produces one outcome. Relocate operator is presented in Figure(1b).

### 3.4.3 Swap

[26] Node Swap is another simple optimisation heuristics: Node swap operator exchanges two locations. Node Swap move is a special case of two subsequent 2-opt moves: the first including both cities and the second without them. It involves removing four links and adding four new links. Therefore it is a specific type of 4-opt move. Swap operator is illustrated in Figure 1c.

### 3.4.4 Three permutations

[22] idea is that three nodes can generate six different orders. Therefore given an existing sequence of three nodes, can create five new solutions. Figure 1d shows that changing the original order of three nodes can result in a better solution.
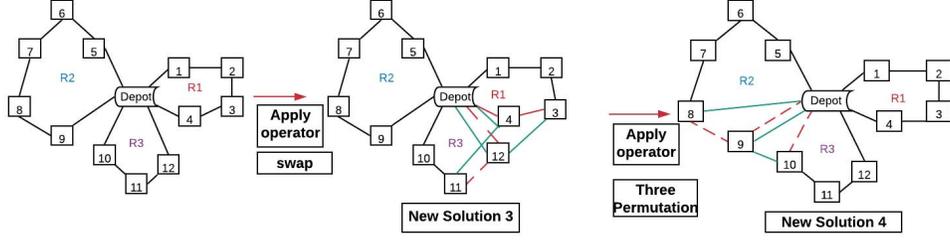
Figure 2: An illustration of a problem instance after applying a set of operators for two routes. Current solution changes to a new solution with dashed lines replaced by red line for the problem CVRP respectively.

## 4 Learning to Guide Local Search (L2GLS)

This section introduces our "Learning to Guide Local Search" (L2GLS) algorithm. Figure 3 illustrates an overview of our approach. Our algorithm starts with a feasible solution and continuously improves the solution using LS operators. L2GLS exploits problem and search-related information to escape local minimas by augmenting the objective function of the problem with a set of penalty terms. The penalty terms are adaptively adjusted and reflects the degree that a solution feature causes a solution to be sub-optimal and candidate for further local improvement. This avoids ad-hoc perturbations to escape local minimal. We found this strategy is able to converge to good solutions faster.

Our method follows threshold-based rules to decide whether we should continuously improve the current solution using LS operators or apply the penalty term. If our RL-based manager decides that the current solution could still be improved, it chooses one LS operator to improve it. When no objective value reduction is made for a pre-determined $I$ number of applied actions (which are operators applied), we consider that the LS has reached a local optimal and the RL-based manager chooses to use penalty term to escape this local optimal. If after M times of reaching a local optimal and applying the penalty terms still result in no cost reduction across these M runs, then the algorithm and search terminates and the best solution of these runs is returned. Having described our overall method, we are now ready to present the details of our RL-based manager.

### 4.1 Reinforcement Learning based Manager

The main component of L2GLS is a RL based manager that guides the search by adaptively selecting which LS operators to use in each iteration and their parameters (e.g, which nodes to apply the operator). It will also detect if there has not been improvement in the solution for a few iterations to determine when we reach a local optima. Then uses the penalty term with adaptive weights to escape the local optima. We will detail these components in the following.

The RL-based manager is a deterministic model, which is defined by the tuple $(S, A, T, R, P)$, where $S$ is the states, $A$ is the actions, $T$ is the deterministic transaction function $(T : S \times A \rightarrow S)$, $R$ is the reward function $(R : S \times A \rightarrow \mathcal{R})$ and $\pi$ is the policy $(\pi : S \rightarrow A)$. We define each of this and relate them to the optimisation problem.

#### 4.1.1 States

Each state represents the problem instance and a solution, and made up of problem-specific and solution-specific features. Problem-specific features include the node location, and demand of each node for CVRP, and they are considered problem-specific because they are associated with the problem and invariant across solutions. Solution-specific features are based on the given solution. For each node, we compute its neighbouring nodes that are visited before and afterwards along with their relevant distances. State also includes recently taken actions and the effects of the actions. We follow the L2I [17] approach for the state. Next, we describe the features. First, we introduce the description of all the features for TSP, where the location of node $i$ is represented by $(m_i, n_i)$, for a 2D location, but this can be generalised to more than 2 dimensional spaces. Previously visited location of node $i$ is denoted by $(m_{i^-}, n_{i^-})$ and the location visited after $i$ is $(m_{i^+}, n_{i^+})$. The neighbouring nodes distance from $i^-$ to $i$ is $(d_{i^-,i})$, distance from $i$ to $i^+$ is $(d_{i,i^+})$, and distance from $i^-$ to $i^+$ is $(d_{i^-,i^+})$. The state also includes the action taken previously and the effects caused by taking the step, i.e., the actions taken $h$ steps before, denoted as $a_{t-h}$, and their effects is denoted by $e_{t-h}$. For example,

6

$a_{t-h}$, $1 \leq h \leq H$, is the action taken h steps before current step t, and its effect $e_{t-h}$ is +1 if the action led to a reduction of total distance, -1 otherwise. All the features for CVRP are the same as TSP, except along with the customer location $(m_i, n_i)$, we have additional tuple dimensions of the demand of each customer $i$, denoted by $G_i$, and the free capacity of the route containing customer $i$ denoted as $C_i$.

### 4.1.2 Actions

L2GLS actions involve selection of appropriate LS operators and their operands to apply to current solution to improve it. Each action is a tuple $(o, \phi)$, where $o \in O$ are the set of LS operators and $\phi \in \Phi$ is the parameters associated with the operators, e.g., remove two edges and reconnect their endpoints in the 2-opt operator (note that $\phi$ will be different for each operator). For a given problem, the same operator with different parameters may perform differently. Consequently, each operator with different parameters act as separate actions and we let the RL model learn how to use them best. Please refer to the Supplementary material to see a description of how the LS operators work.

### 4.1.3 Transition Function and Policy Networks

The transition function is deterministic, i.e., given current state and an action, it will always transit to only one next state. However in terms of learning a policy to guide the actions, the state and action spaces are very large, and it is not feasible to be able to see every possible combination of these to train the policy. Hence we adopt a policy gradient approach to do so, i.e., use a neural network to learn and represent the policy. In a policy gradient approach, a neural network, also known as a policy network, represents the policy, i.e., the mapping function from input states to output actions. This network is defined by a set of hyper-parameters $\theta$ and outputs a probability distribution of actions given an input state. Given a state $S$, our model defines a probability distribution $P_\theta(A|S)$, from which we can sample actions to obtain a solution tour $\pi$. In order to train our policy network, we define the loss of the policy as $\mathcal{L}(\theta|S) = E_{\pi \sim P_\theta(.|S)}[L(\pi|S)] \ \forall a \in A$ (where . is over all actions), which is the expectation of solution ($\pi$) given the probability distribution of $p_\theta(.|S)$ from which we can obtain actions given state. We optimise $\mathcal{L}$ by gradient descent, using the REINFORCE [43] gradient estimator with baseline $b(S)$:

$$\nabla \mathcal{L}(\theta|S) = E_{\pi \sim P_\theta(.|S)}[(L(\pi|S) - b(S)))\nabla_\theta log P_\theta(\pi|S)] \tag{2}$$

We sample solutions and compare it with baseline $b(S)$. We use baseline $b(S)$ to reduce gradient variance (adjusting the probability proportion to result we get compare to the baseline) and speed the learning process. We use ADAM [44] as optimiser. Problem-specific and solution-specific states are transformed into an embedding. The embeddings are then fed to the attention network (see Supplementary for the overview of the policy network). In the policy network, we first fed the problem-specific states to the attention network as an embedding, then the output of the attention network is concatenated with a sequence of recent actions and their effects (solution specific states). Lastly, the output of the attention network is fed into a network of two fully connected layers, producing action probabilities.

### 4.1.4 Rewards

We experiment with two reward functions that were also used in [17]. The first reward is if an selected action (operator and parameters) results in a reduction of the objective value of the current solution, the reward is +1; otherwise, -1. The second reward is defined over a larger range of values. During the first improvement iteration, the objective value attained for the solution is taken as a baseline. For each subsequent iteration, the reward for an operator is equal to the difference between the objective values achieved after applying the operator and the baseline. Similar to [17], we found there is generally a diminishing improvement in latter iterations hence we give the same reward for all operators applied in one iteration.

## 4.2 Penalty Term

The penalty term aims to identify solution features that contribute significantly towards the overall cost of a solution and encourage their replacement in the tour. It does this by adding penalty terms to these features.

The penalty term uses two types of information: one is the cost of each feature to the solution and the set of features currently in the solution. This information is transformed into constraints on features which then are incorporated in the cost function using the adjusted penalty terms. Constraints on features are made possible by augmenting the cost function $L(\pi)$ of the problem to include a set of penalty terms.

The high-cost features are being penalised. For example, in the TSP, the features are the travel costs between pairs of cities, and generally high intercity distances in a tour are not desirable, though we can not exclude them from the beginning because they may be needed to connect remote cities in the tour.

When the search is stuck in a local minimum, our RL manager modifies the parameters on the penalty term [21] to encourage/discourage certain solution features and subsequently calls on LS to improving the solution using the augmented objective function, defined as follows:

$$h(\pi) = L(\pi) + \lambda \cdot \sum_{i=1}^{F} p_i \cdot I_i(\pi) \tag{3}$$

where $F$ is the number of features for the problem instance, $p_i$ is the penalty parameter corresponding to feature $f_i$ and $\lambda$ is a regularisation parameter for the penalty terms. An indicator function $I_i(\pi)$ indicates whether the solution feature $i$ is part of the solution $\pi$ (hence we only penalise and consider features that are part of the current solution).

$$I_i(\pi) = \begin{cases} 1, & \text{solution } \pi \text{ has feature } i \\ 0, & \text{otherwise} \end{cases} \tag{4}$$

The penalty parameter $p_i$ gives the degree up to which the solution feature $f_i$ is constrained. The regularisation parameter $\lambda$ represents the relative importance of penalties with respect to the solution cost and has great significance because it controls the influence of the information on the search process. L2GLS iteratively uses LS to improve the solution and increments some of the parameter parameters $p_1, \cdots, p_F$, each time LS converges to a local minimum. To determine which penalty parameter to increment, the following utility function is computed for each feature:

$$U(\pi_L, f_i) = I_i(\pi_L) \cdot \frac{d_i}{1 + p_i} \tag{5}$$

where $d_i$ is the cost (e.g., intercity travel distance for TSP) of the feature $f_i$. When at a local minimum (denoted as $\pi_L$) recall we seek to penalise those solution features that contribute most towards the current solution. The feature(s) that have the maximal utility have their penalty parameters incremented to penalise them (see Equation 3), which is used by our RL based manager to guide LS operators. But we also do not wish to continually penalise the same features (ones with the largest $d_i$). Hence the utility function is proportional to feature cost $d_i$, but weighted by the inverse of $(1 + p_i)$, which gets larger as we penalise a feature more, and subsequently reduces the utility function as that feature is penalised more. This allows other features to be considered over time to be penalised and increases the searching.

Initially, all the penalty parameters are set to 0 (i.e., no features are constrained), and when LS is searching, a call is made to LS to find a local minimum of the augmented cost function using Equation 3. After each local minimum, the algorithm takes a modified action (with LS operators) on the augmented cost function and uses LS starting from the previously found local minimum. Information is inserted in the augmented cost function by selecting which penalty parameters to increment (e.g., sources of information are the cost of features (edges between two nodes cost) and the local minimum itself). That is how our algorithm offers guidance to the LS.

## 5 Evaluation Setup

For the two problems we are focusing on in this paper, TSP and CVRP, we follow the evaluation setups of AM [11] and L2I [17], which are two state-of-the-art approaches for these problems. All reported metrics, such as the final travelling cost and the running time, are always computed as the average over 1000 randomly generated instances/samples.

### 5.1 Problem Specific Setup

For both problems, we generate the training instances where the coordinates of each node are sampled in a uniformly random way from the unit square $[0; 1] \times [0; 1]$. For TSP, we generated instances with N = 20; 50; 100 nodes and consider these as a small sized TSP problem. Most of the previous work evaluated their models with these sizes, so we seek to first show the performance of L2GLS on small problems. For a large problem, we generated instances with N = 200; 500; 1000 nodes.

For CVRP, we solve the problems with a setup as described in L2I [17]. The location $(m_i; n_i)$ of each customer, as well as of the depot, is uniformly sampled from a unit square (specifically, $m_i; n_i$ are uniformly

distributed in the interval [0;1], respectively), and the travelling cost between two locations $d_{i;j}$ is simply the corresponding distance. The demand $G_i$ of each customer is uniformly sampled from the discrete set $\{1, \cdots, 9\}$. We consider instances with $N = 20; 50; 100$ nodes as a small problem, where the capacity of a vehicle is 20, 30, 40 for $N = 20; 50; 100$, respectively. Similar to TSP, a large problems consists of $N = 200; 500; 1000$ nodes. We keep the capacity of the vehicle 50 for all sizes of the large problems.

### 5.2 Policy Network Settings

Our method is directly comparable to L2I [17], hence to train the model we use the same parameters setting for fair comparison. We use ADAM with a learning rate of 0.001. Unless otherwise stated, we randomly initiate a feasible solution for a problem instance and a given policy and then iteratively update the solution following the policy $M = 40000$ times. After $I = 6$ consecutive improvement steps, we use the penalty term. In After a maximum number of rollout steps, the algorithm stops and among all the 40000 visited solutions, we choose the best as the final solution for a given problem instance. We selected the values of $\lambda$ as constant (in the equation 3), because with the constant value 0.3 was recorded best performance during training. Performance metrics used in the paper includes total travelling cost/distance and the computational running time, which are computed as the average over 1000 random instances. For encouraging exploration, we use a greedy [45] approach, where the RL manager will choose a random improvement action with a probability of 0.05 following [17]. Lastly, L2GLS was implemented in Python, on a workstation with an Intel Xeon 2.4 GHz CPU with 56 cores.
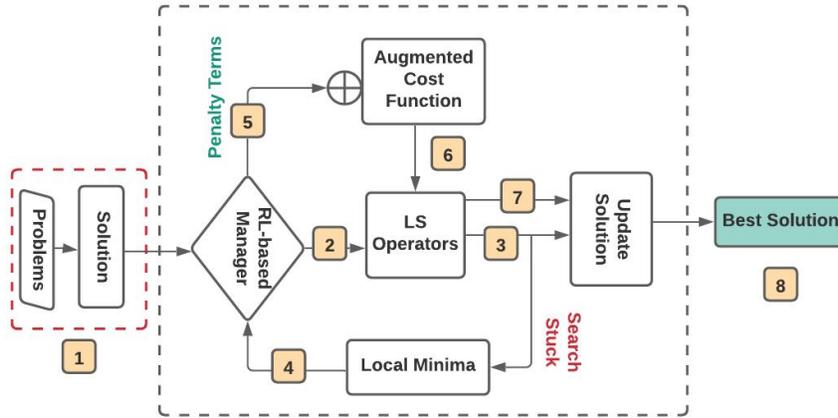


Figure 3: Overview of the Learning to Guide LS Algorithms. L2GLS hierarchy framework. **1:** Given a problem instance, our algorithm first generates a feasible solution. **2:** RL-based manager selects a search operator. **3:** Update the solution **4:** After $I$ improvement steps if no reduce in travelling length means LS stuck in a local minimum. **5:** The RL-based manager again take control to take decision. **6:** RL-based chooses to include a set of penalty terms. **7:** LS is called again to minimise the cost function. **8**: After a certain number of steps (M steps), the model choose the best solutions.

## 6 Experiments and Analysis

The evaluation procedure used in this study involves:

- We evaluated our method using the large instance to show the method's scalability. Few papers considered the large scale instances for TSP and CVRP;

- The model is evaluated using various small scale random instances, as previous neural network-based approaches typically focus on randomly generated small scale data.

- Computational running time analysis for small and large scale instances.

- To demonstrates generalisation ability, we compare the quality of solutions of our method with the benchmark TSP instances from the TSPLIB library [27], where we tested 'thirty-nine' Euclidean

Table 2: Experiment results on small TSP. Average tour length, gap percentage (lower is better, bold is best). Results with * are reported from other papers. (**) means the value reported for the problem size outperformed all the methods, including the optimal one for TSP. (−) means Bresson et al. [34] method never tested on TSP20 instance.

| Method | TSP20 | | TSP50 | | TSP100 | |
|---|---|---|---|---|---|---|
| **Solver** | TourL | Gap | TourL | Gap | TourL | Gap |
| Concorde [7] | 3.84 | 0.00% | 5.70 | 0.00% | 7.77 | 0.00% |
| **Heuristics** | | | | | | |
| LKH3 [8] | 3.84 | 0.00% | 5.70 | 0.00% | 7.77 | 0.00% |
| **Constructive (SL)** | | | | | | |
| GCN* [46] | 3.86 | 0.52% | 5.87 | 2.98% | 8.41 | 8.23% |
| NETSP-Net [31] | 3.85 | 0.26% | 5.85 | 2.63% | 8.31 | 6.94% |
| Fu et al.* [38] | 3.83 | ** | 5.69 | ** | 7.76 | **% |
| **Constructive RL)** | | | | | | |
| AM [11] | 3.87 | 0.78% | 5.80 | 1.75% | 8.15 | 4.89% |
| ERRL [32] | 3.83 | 0.78% | 5.73 | 1.75% | 7.85 | 4.89% |
| POMO[33] | 3.83 | ** | 5.73 | 1.75% | 7.84 | 7.07% |
| Bresson et al.* [34] | - | - | 5.69 | ** | 7.81 | 7.07% |
| **Improvement** | | | | | | |
| LHI*[16] | 3.83 | ** | 5.74 | 0.71% | 8.01 | 3.08% |
| L2I. [17] | 3.84 | 0.00 | 5.72 | 0.7% | 7.90 | 1.67% |
| RL2OPT.*[35] | 3.84 | 0.00 | 5.72 | 0.7% | 7.91 | 1.80% |
| **L2GLS(Ours)** | **3.72** | ** | **5.65** | ** | **7.69** | ** |

instances. For the first time, in the L2O field, we evaluated our model up to 1002 benchmark instances, achieving close to the optimal solution. Moreover, we also compare the solution quality of our method using various scenarios of CVRP benchmark instances from [28]. Demonstrate the model generalisation using various training and testing sizes of instances for TSP and CVRP. Due to space limitations, we presented these results in a supplementary document.

- We illustrate the impact of L2GLS variants and evaluated on randomly generated instances for TSP and CVRP. Moreover, we compare L2GLS performance with and without using penalty terms on TSPLIB [27].

Table 3: Experiment results on small CVRP. Average tour length and gap percentage (lower is better, best in bold). Results with * are reported from AM ([11]) and Chen et al. [18]. (**) means the value reported for the problem size outperformed all the methods, including the optimal one for CVRP.

| Method | CVRP20 | | CVRP50 | | CVRP100 | |
|---|---|---|---|---|---|---|
| **Solver** | TourL | Gap(%) | TourL | Gap(%) | TourL | Gap(%) |
| LKH3 [8] | 6.14 | 0.00 | 10.39 | 0.00% | 15.67 | 0.00 |
| **Heuristics** | | | | | | |
| Or-tools* | 6.43 | 4.73 | 11.43 | 10.00 | 17.16 | 9.50 |
| **Constructive** | | | | | | |
| AM [11] | 6.67 | 8.63 | 11.00 | 5.87 | 16.99 | 8.42 |
| Nazari et al. [13] | 7.07 | 15.14 | 11.95 | 15.01 | 17.89 | 14.16 |
| POMO[33] | 6.35 | 3.42 | 10.74 | 3.36 | 16.15 | 3.06 |
| ERRL [32] | 6.34 | 3.25 | 10.77 | 3.65 | 16.35 | 4.02 |
| **Improvement** | | | | | | |
| Chen et al. [18] | 6.16 | 0.3 | 10.51 | 1.15 | 16.10 | 2.72 |
| LIH. [16] | 6.16 | 0.26 | 10.72 | 0.35 | 16.30 | 1.60 |
| L2I [17] | 6.12 | ** | 10.35 | ** | 15.57 | ** |
| **L2GLS (Ours)** | **5.85** | ** | **10.30** | ** | **14.67** | ** |

Table 4: Experiment results of average tour length on large TSP and CVRP (lower is better, best in bold).

(a) TSP

| Method | TSP200 | TSP500 | TSP1000 |
|---|---|---|---|
| **Solver** | TourL | TourL | TourL |
| Concorde [7] | 10.15 | 16.542 | 23.130 |
| L2I [17] | 10.95 | 17.68 | 26.13 |
| **L2GLS** | **10.10** | **16.12** | **22.37** |

(b) CVRP

| Method | CVRP200 | CVRP500 | CVRP1000 |
|---|---|---|---|
| **Solver** | TourL | TourL | TourL |
| L2I [17] | 28.4 | 64.68 | 128.49 |
| **L2GLS** | **24.69** | **60.67** | **124.29** |

## 6.1 Performance analysis: Small scale TSP and CVRP instances

In this section, we report the performance of our method on various sizes of small scale instances for TSP and CVRP.

(a) TSP computation time.



(b) CVRP computation time.

Figure 4: Average computation time for smaller TSP and CVRP.

### 6.1.1 TSP

Table 2 reports the performance of L2GLS and baseline methods. We have five different baseline approaches in Table 2: solver Concorde [7], heuristics [8], learning methods using constructive heuristics, supervised and RL approaches; and improvement heuristics namely LHI [16], RL2OPT [35] and L2I [17]. There are several ML models in the literature, but they all produce poor performance compare to our approach, thus being omitted here. Our approach outperformed all recent learning method, presented in Table 2. Note that the L2GLS method produces state-of-the-art results for TSP and to the best of our knowledge is the first learning-based framework that outperforms Concorde [10] for the random dataset.

### 6.1.2 CVRP

Table 3 reported the performance of L2GLS methods for CVRP, there we compare our method against the current RL algorithms [13] [11] [17] [16] [32]. In particular, the average tour length achieved by L2GLS is substantially shorter than LKH3 [8]. L2GLS also outperforms all the neural network-based approaches, including recent works by Chen et al. [18] and L2I [17]. Our algorithm generates state-of-the-art results for TSP and CVRP for small problems.
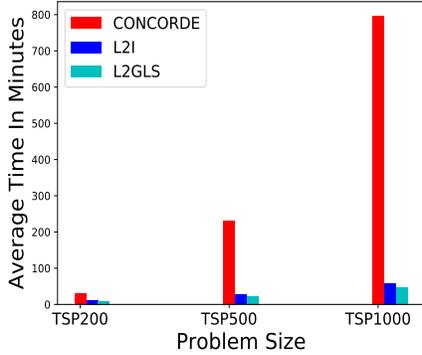
### 6.2 Scalability Analysis: Large Scale TSP and CVRP instances

In this section, we compared the L2GLS performance to a most recent state of the art learning-based approach L2I [17] for large sized problems, as we discussed previously, most state of the art works have not evaluated on due to their lack of scalability. Tables 4a and 4b reports the results of large TSP and CVRP instances respectively. Our method significantly outperforms L2I [17] for all the TSP and CVRP instances, even for a large number of nodes N=1000. For CVRP1000, we averaged over 200 instances instead of 1000 instances due to time constraint.
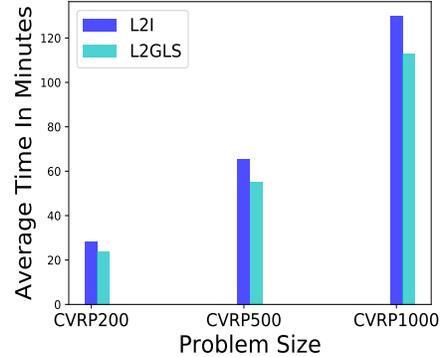
### 6.3 Computational Running Time Analysis

Comparing the algorithms' efficiency in terms of execution time is difficult due to varying hardware and implementations among different approaches. So for fair comparison, we performed experiment on the same hardware we used for our work. We compare the inference time with a recant improvement based approach L2I [17] using their publicly available code [1]. For TSP Figures 4a 5a illustrate that L2GLS is efficient compared to Concorde [10]. For example, to solve TSP100, Concorde [10] needs approximately 60 minutes, whereas L2GLS required only 11 mins to solve (Figure 4a). Given 1000 random instances. L2GLS significantly reduces the running time for all instances compare to Concorde [10] L2I [17] for all problem sizes of TSP without sacrificing finding optimal tours, as shown in Figure 4a 5a. For CVRP, the running time for LKH3 is rapidly increases when the problem size is large (Figure 4b). For example, to solve CVRP100, LKH3 needs 780 minutes, whereas L2GLS requires only 17 minutes, as shown in Figure 4b. L2GLS allows notable reduction in the run time with good quality solution compare to L2I [17].

---

[1] https://github.com/rlopt/l2i

(a) TSP computation time.



(b) CVRP computation time.

Figure 5: Average computation time for larger TSP and CVRP.

We also run experiments with larger instances with sizes up to 1000 nodes where results are summarised in Figure 5a. The Figure 5a shows that L2GLS saves more running time compared to Concorde [7]. For example, in terms of TSP1000 problems, Concorde [10] needs on average 796 minutes to solve a problem. However, L2GLS only need 47 minutes to solve the same problems, as shown in Figure 5a. Compared with L2I [17], L2GLS uses less average running time and tour length.

In the case of CVRP, It is evident from Figure 4b that the running time of LKH3 grows with an increase in number of nodes, i.g., for CVRP100 size, LKH3 takes 780 minutes. Therefore, we fairly compared the running times with the L2I [17] methods to show the efficiency of our method for CVRP200, CVRP500 and CVRP1000. The results are presented in Figure 5b. Specifically, in terms of running time, our method not only is more efficient when compared to LKH3, but also outperformed L2I [17], as illustrated in Figure 4 and 5.

It is worth noting that for smaller and larger randomly generated instances, L2GLS can generate solutions faster than the state of the art algorithms, i.e., Concorde for TSP and LKH3 for CVRP (4 and 5).

## 6.4 Generalisation - L2GLS method

In this section, we will show how our method can generalise well to benchmark libraries. Specifically we utilise TSPLIB instances [27] which are generally used to compare TSP algorithms and data proposed by Uchoa et al. [28] for CVRP. In addition, we evaluate our method's generalisation performance on different problem sizes of TSP and CVRP, where we trained and tested on different size of benchmark datasets.

### 6.4.1 Benchmark (TSPLIB)

L2GLS model was trained on randomly generated 50 nodes of instances. Our aim in these experiments are to know to what extent the learned algorithm generalises to benchmark instances. The distribution of the instance, such as node locations, is completely different from those we used in the training instance. We compare our method with S2VDQN [47], AM [11], LHI [16] and L2I [17], the results are shown in Table 5, which demonstrates that L2GLS outperformed all the learning-based approaches most of the instances. Moreover, the heuristic solver OR-Tools generates inferior solutions compared to L2GLS [27]. We tested L2GLS methods on the 39 TSP instances up to size 1002 instance from TSPLIB [27] in Table 5. For all the scenarios, L2GLS obtain a small gap compare to Concorde [7] and outperformed all the neural network-based approaches. All the recent works only evaluated their methods up to 299 nodes, but to show the capability of our approach, we tested another three instances, namely Lin318, Pr439 and Pr1002. As shows in Table 5, our algorithm outperforms the prior ML based approaches in terms of solution quality. For three large instances we compare our method with L2I [17]. Our L2GLS optimality gap is 2.74% 3.77% and 1.19%, whereas L2I [17] gap is 5.04% 5.50% and 5.76% for Lin318, Pr439 and Pr1002 problem instances respectively. Table 5 also shows L2GLS outperformed Concorde [7] for some instances (denoted as ∗∗ in Table 5), such as pr107, pr144 and u159. However, we emphasise that the aim is not to outperform Concord but to show L2GLS performs well.

Table 5: TSPLIB results: Instances are sorted by increasing size, with the number at the end of an instance's name indicating its size. Values reported are the cost of the tour found by each method (lower is better, best in bold). Gap (%) is the gap to the solution obtained by Concorde [7]. The results marked * reported from S2VDQN [47] and marked † reported from LHI [16]. ∗∗ means (column: Gap%) our model performed better than solver, no gap reported.

| Problems | L2I | | Or-tools | | S2VDQN* | | AM† | | LHI† | | L2GLS | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | TourL | Gap(%) | TourL | Gap(%) | TourL | Gap(%) | TourL | Gap(%) | TourL | Gap(%) | TourL | Gap(%) |
| Eil51 | 430 | 0.93 | 436 | 2.34 | 439 | 3.05 | 435 | 2.11 | 438 | 2.81 | **428** | 0.46 |
| Berlin52 | 7,694 | 2.01 | 7,945 | 5.34 | **7,542** | 0.00 | 7,668 | 1.67 | 8,020 | 6.33 | 7,544 | 0.026 |
| St70 | 683 | 1.85 | 683 | 1.18 | 696 | 3.11 | 690 | 2.22 | 706 | 4.59 | **679** | 0.29 |
| Eil76 | 551 | 2.41 | 561 | 4.27 | 564 | 4.83 | 563 | 4.64 | 575 | 6.87 | **546** | 1.48 |
| Pr76 | 108,871 | 0.65 | 111,104 | 2.72 | 108,446 | 0.26 | 111,250 | 2.85 | 109,668 | 1.39 | **108,159** | 0.00 |
| rat99 | 1,257 | 3.79 | 1,232 | 1.73 | 1,280 | 5.69 | 1,319 | 8.91 | 1,419 | 17.17 | **1,221** | 0.82 |
| KroA100 | 22,036 | 3.54 | **21,448** | 0.78 | 21,897 | 2.88 | 38,200 | 79.49 | 25,196 | 18.39 | 21,650 | 1.72 |
| KroB100 | 22,452 | 1.40 | 23,006 | 3.90 | 22,692 | 2.48 | 35,511 | 60.38 | 26,563 | 19.97 | **22,191** | 0.22 |
| KroC100 | 21,258 | 2.45 | 21,583 | 4.01 | 21,074 | 1.56 | 30,642 | 47.67 | 25,343 | 22.14 | **20,820** | 0.34 |
| KroD100 | 22,102 | 3.79 | 21,636 | 1.60 | 22,102 | 3.79 | 32,211 | 51.26 | 24,771 | 16.32 | **21,514** | 1.03 |
| KroE100 | 22,875 | 3.65 | 22,598 | 2.40 | 22,913 | 3.82 | 27,164 | 23.09 | 26,903 | 21.90 | **22,152** | 4.02 |
| rd100 | 8,132 | 2.80 | 8,189 | 3.52 | 8,159 | 3.14 | 8,152 | 3.05 | **7,915** | 0.06 | 7944 | 0.42 |
| eil101 | 654 | 3.97 | 664 | 5.56 | 659 | 4.76 | 667 | 6.04 | 658 | 4.61 | **642** | 2.06 |
| Lin105 | 14,700 | 2.23 | 14,824 | 3.09 | 15,023 | 4.47 | 51,325 | 256.94 | 18,194 | 26.53 | **14,406** | 0.18 |
| pr107 | 44,807 | 1.13 | 45,072 | 1.73 | 45,113 | 1.82 | 205,519 | 363.89 | 53,056 | 19.75 | **44,301** | ** |
| pr124 | 59,935 | 1.53 | 62,519 | 5.91 | 61,623 | 4.39 | 167,494 | 183.74 | 66,010 | 11.82 | **59,030** | 0.00 |
| bier127 | 120,699 | 2.04 | 122,733 | 3.76 | 121,576 | 2.78 | 207,600 | 75.51 | 142,707 | 20.64 | **119,281** | 0.84 |
| Ch130 | 6,318 | 3.40 | 6,470 | 5.89 | 6,270 | 2.61 | 6,316 | 3.37 | 7,120 | 16.53 | **6.178** | 1.11 |
| pr136 | 98,728 | 2.02 | 102,213 | 5.62 | 99,474 | 2.79 | 102,877 | 6.36 | 105,618 | 9.14 | **98680** | 0.50 |
| pr144 | 59,895 | 2.31 | 59,286 | 1.27 | 59,436 | 1.53 | 183,583 | 213.61 | 71,006 | 21.30 | **58,535** | ** |
| Ch150 | 6,779 | 3.84 | 7,232 | 10.78 | 6,985 | 7.00 | 6,877 | 5.34 | 7,916 | 21.26 | **6,606** | 1.19 |
| KroA150 | 27,307 | 2.95 | 27,592 | 4.02 | 27,888 | 5.14 | 42,335 | 59.61 | 31,244 | 17.79 | **27,248** | 2.72 |
| KroB150 | 26,563 | 19.97 | 23,006 | 3.90 | 27,209 | 22.88 | 35,511 | 60.38 | 26,563 | 19.97 | **22,381** | 1.08 |
| pr152 | 75,196 | 1.05 | 75,834 | 2.92 | 75,283 | 2.17 | 103,110 | 39.93 | 85,616 | 16.19 | **74,204** | 0.70 |
| u159 | 43,220 | 2.70 | 45,778 | 8.78 | 45433 | 7.96 | 115,372 | 174.17 | 51,327 | 21.97 | **42,075** | ** |
| rat195 | 2,403 | 3.44 | 2,389 | 2.89 | 2,581 | 11.10 | 3,661 | 57.59 | 2,913 | 25.39 | **2,368** | 1.93 |
| d198 | 16,349 | 3.60 | 15,963 | 1.15 | 16,453 | 4.26 | 68,104 | 331.58 | 17,962 | 13.82 | **16,244** | 2.94 |
| KroA200 | 30,617 | 4.25 | 29,741 | 1.27 | 30,965 | 5.43 | 58,643 | 99.68 | 35,958 | 22.43 | **29,736** | 1.25 |
| KroB200 | 30,925 | 5.05 | 30,516 | 3.66 | 31,692 | 7.66 | 50,867 | 72.79 | 36,412 | 23.69 | **30,154** | 2.43 |
| ts225 | 129,638 | 2.36 | 128,564 | 1.51 | 136,302 | 7.62 | 141,628 | 11.83 | 158,748 | 25.35 | **126,645** | 0.001 |
| Tsp225 | 3,995 | 2.01 | 4,046 | 3.31 | 4,154 | 6.07 | 24,816 | 533.7 | 4,701 | 20.04 | **3,940** | 0.16 |
| pr226 | 87,160 | 8.44 | 82,968 | 3.23 | 81,873 | 1.87 | 101,992 | 26.90 | 97,348 | 21.12 | **81,756** | 1.72 |
| gil262 | 2,534 | 6.56 | 2,519 | 5.92 | 2,537 | 6.68 | 2,693 | 13.24 | 2,963 | 24.60 | **2,451** | 3.06 |
| pr264 | 56,191 | 14.36 | 51,954 | 5.73 | 52,364 | 6.57 | 338,506 | 588.93 | 65,946 | 34.21 | **50,528** | 0.00 |
| A280 | 2,751 | 6.66 | 2713 | 5.19 | 2,867 | 11.16 | 11,810 | 357.92 | 2,989 | 15.89 | **2,672** | 3.60 |
| Pr299 | 53,753 | 11.54 | 49,447 | 2.60 | **51,895** | 4.95 | 513,673 | 938.83 | 59,786 | 20.90 | 52,827 | 7.68 |
| Lin318 | 44,151 | 5.04 | - | | 45,375 | 7.96 | - | - | - | - | **43,184** | 2.74 |
| Pr439 | 113,124 | 5.50 | - | - | - | - | - | - | - | - | **111,269** | 3.77 |
| Pr1002 | 273,970 | 5.76 | - | - | - | - | - | - | - | - | **262,150** | 1.19 |

### 6.4.2 Practical Instances (CVRP)

This section evaluated L2GLS methods with the CVRP benchmarks [28], where each instance is characterised by the following attributes: number of customers, depot positioning, and customer positioning. The details of each attribute are described in the supplementary. This data distribution is commonly used in the OR field. We have evaluated L2GLS on various scenarios of CVRP benchmark instances. We utilised three sizes of the CVRP dataset with some additional scenarios following the data generation proposed in [28]. We further combined the customer and depot positioning and reported nine scenarios, namely random, central and eccentric depot positioning, combined with random, clustered and random-clustered customer positioning. Table 6 shows the impact of depot positioning while using different customer position. This experiment which demonstrates that L2GLS works over different data distribution ( on nine different data distributions).

### 6.5 Ablations

To evaluate the value of a different set of LS operators, we run four variants of L2GLS, i.e., L2GLS, L2GLS2, L2GLS3 and L2GLS without penalty terms. **L2GLS** consists of ($2opt + relocate + swap + threepermutation$) with penalty term that our main method, **L2GLS2** consists of ($2opt + relocate + swap$), **L2GLS3** consists of ($2opt + swap + threepermutation$), and **L2GLS without penalty terms** consists of operators ($2opt + relocate + swap + threepermutation$) without the penalty term.

Table 6: CVRP instances analysis using nine different scenarios, namely random, central and eccentric depot positioning combine with random, clustered and random-clustered customer positioning. Values reported are the cost of the tour found.

| Problems | => | CVRP20 | CVRP50 | CVRP100 |
|---|---|---|---|---|
| **Depot Position** | **Customer Position** | | | |
| Random | Random | 5.86 | 10.37 | 15.60 |
| Random | Clustered | 4.82 | 8.46 | 13.45 |
| Random | Random-Clustered | 5.86 | 10.31 | 16.59 |
| Central | Random | 5.07 | 8.79 | 13.59 |
| Central | Clustered | 3.74 | 6.43 | 9.96 |
| Central | Random-Clustered | 5.07 | 8.25 | 12.4 |
| Eccentric | Random | 7.42 | 13.53 | 20.89 |
| Eccentric | Clustered | 6.05 | 10.86 | 17.46 |
| Eccentric | Random-Clustered | 7.35 | 13.14 | 19.55 |

Table 7: Comparison of our different methods (combination of LS operators), L2GLS, L2GLS2,L2GLS3, L2GLS without penalty term.

(a) TSP results.

| **TSP** | TourL TSP20 | TourL TSP50 | TourL TSP100 |
|---|---|---|---|
| L2GLS | 3.72 | 5.67 | 7.69 |
| L2GLS2 | 3.80 | 5.72 | 7.80 |
| L2GLS3 | 3.80 | 5.74 | 7.80 |
| L2GLS without penalty | 3.72 | 5.69 | 7.72 |

(b) CVRP results.

| **CVRP** | TourL CVRP20 | TourL CVRP50 | TourL CVRP100 |
|---|---|---|---|
| L2GLS | 5.85 | 10.30 | 14.67 |
| L2GLS2 | 6.12 | 10.5 0 | 16.54 |
| L2GLS3 | 6.17 | 11.00 | 17.70 |
| L2GLS without penalty | 5.86 | 10.37 | 15.60 |

The goal of this experiment is to analyse the impact of a different set of LS operators. As we mentioned earlier we implemented and experimented with a set of LS operators and concluded with the L2GLS method, which is a combination of 2-opt, relocate, swap, and three permutations with penalty terms that is computationally efficient and able to produce state-of-the-art results for TSP.

### 6.5.1 Evaluation of the four variants of L2GLS on Random TSP and CVRP instances

Tables 7a and 7b illustrate the performance of a different set of search operators to solve random TSP and CVRP problem instances. Each variant used the same neural network architecture, to evaluate the effect of the search operators on the solving results. It can be observed that using a different variant of operators impacts the solution quality. The best TSP and CVRP results are obtained by L2GLS among all the variants. From Tables 7a and 7b, it can be seen that L2GLS3, by not using relocate, generates significantly worsen results for TSP and CVRP.

### 6.5.2 Evaluation of two variants of L2GLS on TSPLIB instances

Among the four variants, L2GLS and L2GLS without the penalty term performed close to the optimal solution Concorde [7] shown in Tables 7a and 7b for random instances. Therefore, we selected benchmark from our Table 5 and evaluated our model without penalty term to check how our model performed. From Table 8, it is clear that using penalty term with LS operators significantly improves over only using LS operators to solve the problems.

We demonstrated that our method outperformed recent approaches. The potential for better performance is the choice of heuristics used as LS operators. As many LS operators are effective and efficient for solving COP, therefore, in [25] states that well-implemented LS creation can compete with the best heuristics. Misztal et al., [20] studied that swap belongs to a group of the LS algorithm characterised by the good quality of the returned solution. In [22] they concluded the 2-opt is the best individual operator, besides the mixed variant of operators found the optimum solution for many cases. In this work, we investigate different operators, which is a mixed variant of operators ($2opt + relocate + swap + threepermutation$). Moreover, we have introduced penalty terms, which contributes to significantly improved results for RPs. When searching stuck in local minima, we use penalty terms to improve the solution further; therefore, our result shows that these combined operators yield better solutions for most instances from TSPLIB [27].

Table 8: Compared two L2GLS variants: L2GLS and L2GLS without penalty term, results are reported for TSPLIB. ∗∗ means (column: Gap%) our model performed better than solver, no gap reported. Bold means better.

| Problems | Concorde | | L2GLS without Penalty | | L2GLS | |
|---|---|---|---|---|---|---|
| | TourL | Gap(%) | TourL | Gap(%) | TourL | Gap(%) |
| St70 | 675 | 0.00 | 679 | 0.59 | **677** | 0.29 |
| Pr76 | 108,159 | 0.00 | 108,570 | 0.37 | **108,159** | 0.00 |
| rat99 | 1,211 | 0.00 | 1,223 | 0.99 | **1,221** | 0.82 |
| KroA100 | 21,282 | 0.00 | 21,866 | 2.74 | **21,650** | 1.72 |
| rd100 | 7,910 | 0.00 | 8,010 | 1.26 | **7944** | 0.42 |
| Lin105 | 14,379 | 0.00 | 14,565 | 1.29 | **14,406** | 0.18 |
| pr107 | 44,303 | 0.00 | 44,527 | 0.50 | **44,301** | ** |
| Ch130 | 6,110 | 0.00 | 6,192 | 1.39 | **6,178** | 1.11 |
| pr144 | 58,537 | 0.00 | 58,781 | 0.44 | **58,535** | ** |
| KroA150 | 26,524 | 0.00 | 27,408 | 3.33 | **27,248** | 2.72 |
| u159 | 42,080 | 0.00 | 42,645 | 1.35 | **42,075** | ** |
| d198 | 15,780 | 0.00 | 16,341 | 3.55 | **16,244** | 2.94 |
| ts225 | 126,643 | 0.00 | 129,248 | 2.05 | **126,645** | 0.001 |
| Tsp225 | 3,916 | 0.00 | 3,940 | 2.83 | **3,940** | 0.16 |
| Pr299 | 48,191 | 0.00 | 52,827 | 9.62 | 51,895 | 7.68 |

Table 9: Trained with TSP50. Average tour length, gap percentage (lower is better, best in bold).

| | TourL | TourL | TourL |
|---|---|---|---|
| **TSP** | TSP20 | TSP50 | TSP100 |
| Concorde | 3.84 | 5.70 | 7.77 |
| L2GLS | 3.72 | 5.69 | 7.72 |
| **CVRP** | CVRP20 | CVRP50 | CVRP100 |
| LKH3 [8] | 6.14 | 10.39 | 15.67 |
| L2GLS | 5.86 | 10.52 | 15.9 |

## 6.6 Various test and training

One of the aims of the proposed methods was to ensure its generalisation on different problem sizes of TSP and CVRP. We trained the model with TSP50, used to solve TSP20, TSP50, TSP100, and CVRP20, CVRP50, CVRP100 presents the result in Table9. We trained the same model with TSP50, used to solve TSP20, TSP50, TSP100, CVRP20, CVRP50, CVRP100, and present the result in Table10. We show that L2GLS can generalise to different problem distributions even when trained on different sizes and problems.

## 7 Conclusion

We proposed L2GLS for solving RPs, which starts with an initial solution and improves the solution with LS operators selected by an RL-based manager or with penalty terms to guide the LS further to improve the solution. We also analysed distinct choices of LS operators and selected the best choice of operators from this analysis. The current need is to find an effective combination of LS operators that can generate optimum solutions that adapt and generalise to different problems, which our reinforcement-based manager achieved. Among the four design choices discussed in Section 6.5, we selected L2GLS because we achieve better solution quality applied L2GLS on TSP and relative better for CVRP. Many learning algorithms can solve small, randomly generated problems. However, many applications in the real world are related to larger instances; there is also variability in data distribution. Most of the state of the art algorithms have challenges to solve these practical problems. Therefore, we proposed a method to choose effective LS operators and let RL-based managers learn all the problems to predict various other RPs. L2GLS achieved new state-of-the-art results for TSP instances. Moreover, it outperforms not only for small scale problems but can generalise well benchmark and large-scale datasets. For CVRP, it outperforms LKH3, OR-tools [48] and recent DL-based baselines.

Future research can explore expanding our solution framework to solve other variants of combinatorial problems. Moreover, advanced search operators can be utilised for other combinatorial problems.

## References

[1] Roberto Baldacci, Enrico Bartolini, and Gilbert Laporte. Some applications of the generalized vehicle routing problem. *Journal of the operational research society*, 61(7):1072–1077, 2010.

Table 10: Trained with CVRP50.Average tour length, gap percentage (lower is better, best in bold)

|  | TourL | TourL | TourL |
|---|---|---|---|
| **TSP** | TSP20 | TSP50 | TSP100 |
| Concorde | 3.84 | 5.70 | 7.77 |
| L2GLS | 3.76 | 5.70 | 7.76 |
| **CVRP** | CVRP20 | CVRP50 | CVRP100 |
| LKH3 | 6.14 | 10.39 | 15.67 |
| L2GLS | 5.85 | 10.30 | 14.67 |

[2] George Dantzig, Ray Fulkerson, and Selmer Johnson. Solution of a large-scale traveling-salesman problem. *Journal of the operations research society of America*, 2(4):393–410, 1954.

[3] Michael Held and Richard M Karp. A dynamic programming approach to sequencing problems. *Journal of the Society for Industrial and Applied mathematics*, 10(1):196–210, 1962.

[4] Georges A Croes. A method for solving traveling-salesman problems. *Operations research*, 6(6):791–812, 1958.

[5] Marshall L Fisher. The lagrangian relaxation method for solving integer programming problems. *Management science*, 27(1):1–18, 1981.

[6] Fred Glover and Eric Taillard. A user's guide to tabu search. *Annals of operations research*, 41(1):1–28, 1993.

[7] David Applegate, Ribert Bixby, Vasek Chvatal, and William Cook. Concorde tsp solver, 2006.

[8] Keld Helsgaun. An extension of the lin-kernighan-helsgaun tsp solver for constrained traveling salesman and vehicle routing problems. *Roskilde: Roskilde University*, 2017.

[9] Shen Lin and Brian W Kernighan. An effective heuristic algorithm for the traveling-salesman problem. *Operations research*, 21(2):498–516, 1973.

[10] David L Applegate, Robert E Bixby, Vasek Chvatal, and William J Cook. *The traveling salesman problem: a computational study*. Princeton university press, 2006.

[11] WWM Kool and M Welling. Attention solves your tsp. *arXiv preprint arXiv:1803.08475*, 2018.

[12] Irwan Bello, Hieu Pham, Quoc V Le, Mohammad Norouzi, and Samy Bengio. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.

[13] Mohammadreza Nazari, Afshin Oroojlooy, Lawrence Snyder, and Martin Takác. Reinforcement learning for solving the vehicle routing problem. In *Advances in Neural Information Processing Systems*, pages 9861–9871, 2018.

[14] David SW Lai, Ozgun Caliskan Demirag, and Janny MY Leung. A tabu search heuristic for the heterogeneous vehicle routing problem on a multigraph. *Transportation Research Part E: Logistics and Transportation Review*, 86:32–52, 2016.

[15] Lijun Wei, Zhenzhen Zhang, Defu Zhang, and Stephen CH Leung. A simulated annealing algorithm for the capacitated vehicle routing problem with two-dimensional loading constraints. *European Journal of Operational Research*, 265(3):843–859, 2018.

[16] Yaoxin Wu, Wen Song, Zhiguang Cao, Jie Zhang, and Andrew Lim. Learning improvement heuristics for solving routing problems. *arXiv preprint arXiv:1912.05784*, 2019.

[17] Hao Lu, Xingwen Zhang, and Shuang Yang. A learning-based iterative method for solving vehicle routing problems. In *International Conference on Learning Representations*, 2019.

[18] Xinyun Chen and Yuandong Tian. Learning to perform local rewriting for combinatorial optimization. In *Advances in Neural Information Processing Systems*, pages 6281–6292, 2019.

[19] André Hottung and Kevin Tierney. Neural large neighborhood search for the capacitated vehicle routing problem. *arXiv preprint arXiv:1911.09539*, 2019.

[20] Wojciech Misztal. The impact of perturbation mechanisms on the operation of the swap heuristic. *Archiwum Motoryzacji*, 86(4), 2019.

[21] Christos Voudouris and Edward Tsang. Guided local search and its application to the traveling salesman problem. *European journal of operational research*, 113(2):469–499, 1999.

[22] Lahari Sengupta, Radu Mariescu-Istodor, and Pasi Fränti. Which local search operator works best for the open-loop tsp? *Applied Sciences*, 9(19):3985, 2019.

[23] Helena Ramalhinho Lourenço, Olivier C Martin, and Thomas Stützle. Iterated local search: Framework and applications. In *Handbook of metaheuristics*, pages 129–168. Springer, 2019.

[24] Michel Gendreau, Alain Hertz, and Gilbert Laporte. New insertion and postoptimization procedures for the traveling salesman problem. *Operations Research*, 40(6):1086–1094, 1992.

[25] Florian Arnold and Kenneth Sörensen. Knowledge-guided local search for the vehicle routing problem. *Computers & Operations Research*, 105:32–46, 2019.

[26] Mauricio GC Resende and Renato F Werneck. A fast swap-based local search procedure for location problems. *Annals of Operations Research*, 150(1):205–230, 2007.

[27] Gerhard Reinelt. Tsplib—a traveling salesman problem library. *ORSA journal on computing*, 3(4):376–384, 1991.

[28] Eduardo Uchoa, Diego Pecin, Artur Pessoa, Marcus Poggi, Thibaut Vidal, and Anand Subramanian. New benchmark instances for the capacitated vehicle routing problem. *European Journal of Operational Research*, 257(3):845–858, 2017.

[29] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Advances in Neural Information Processing Systems*, pages 2692–2700, 2015.

[30] 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.

[31] Nasrin Sultana, Jeffrey Chan, and AK Qin. Learning to optimise general tsp instances. *arXiv preprint arXiv:2010.12214*, 2020.

[32] Nasrin Sultana, Jeffrey Chan, AK Qin, and Tabinda Sarwar. Learning vehicle routing problems using policy optimisation. *arXiv preprint arXiv:2012.13269*, 2020.

[33] Yeong-Dae Kwon, Jinho Choo, Byoungjip Kim, Iljoo Yoon, Seungjai Min, and Youngjune Gwon. Pomo: Policy optimization with multiple optima for reinforcement learning. *arXiv preprint arXiv:2010.16011*, 2020.

[34] Xavier Bresson and Thomas Laurent. The transformer network for the traveling salesman problem. *arXiv preprint arXiv:2103.03012*, 2021.

[35] Paulo R de O da Costa, Jason Rhuggenaath, Yingqian Zhang, and Alp Akcay. Learning 2-opt heuristics for the traveling salesman problem via deep reinforcement learning. *arXiv preprint arXiv:2004.01608*, 2020.

[36] 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.

[37] Ke Li and Jitendra Malik. Learning to optimize. *arXiv preprint arXiv:1606.01885*, 2016.

[38] 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.

[39] Wouter Kool, Herke van Hoof, Joaquim Gromicho, and Max Welling. Deep policy dynamic programming for vehicle routing problems. *arXiv preprint arXiv:2102.11756*, 2021.

[40] Pierre Hansen and Nenad Mladenović. First vs. best improvement: An empirical study. *Discrete Applied Mathematics*, 154(5):802–817, 2006.

[41] Ibrahim H Osman and James P Kelly. Meta-heuristics: an overview. *Meta-heuristics*, pages 1–21, 1996.

[42] Michel Gendreau and Jean-Yves Potvin. Metaheuristics in combinatorial optimization. *Annals of Operations Research*, 140(1):189–213, 2005.

[43] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.

[44] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[45] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[46] Chaitanya K Joshi, Thomas Laurent, and Xavier Bresson. An efficient graph convolutional network technique for the travelling salesman problem. *arXiv preprint arXiv:1906.01227*, 2019.

[47] Elias Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, and Le Song. Learning combinatorial optimization algorithms over graphs. In *Advances in Neural Information Processing Systems*, pages 6348–6358, 2017.

[48] Laurent Perron and Vincent Furnon. Google's or-tools. *URL https://developers. google. com/optimization*, 2019.

[49] Marius M Solomon. Algorithms for the vehicle routing and scheduling problems with time window constraints. *Operations research*, 35(2):254–265, 1987.
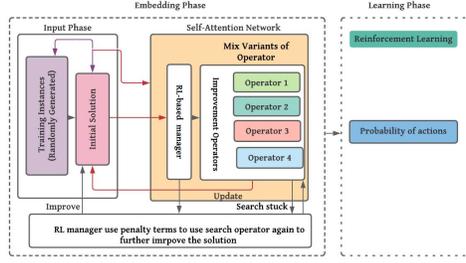
Figure 6: Overview of policy network. The first box is the state embedding part of policy network. Embedding part contains problem features and solution features, an attention network, and a sequence of actions and their effects.

## .1 Policy Network

In our policy network problem and solution-specific input features are transformed into an embedding, which is fed into an attention network. The output of the attention network is concatenated with a sequence of recent actions and their effects. In the end, the output of the attention network is fed into two fully connected layers. Figure 6 shows our overall policy network.

## .2 Instance attributes: CVRP

Instance attributes, depot positioning Three different positions for the depot are considered: Central (C) – depot in the centre of the grid, Eccentric (E) – depot in the corner of the grid, point and Random (R) – depot in a random point of the grid. Customer positioning Three alternatives for customer positioning are considered: following the R, C and RC instance classes of the Solomon set for the VRPTW [49]. Random (R) – all customers are positioned at random points of the grid. Clustered (C) – A number N of customers that will act as cluster seeds are picked from a uniform discrete distribution. Next, the N nodes are randomly positioned in the grid following [28]. Random-clustered (RC)– Some customers are clustered using Cluster positioning; the remaining customers are randomly positioned.