# Deep Reinforcement Learning for Solving the Vehicle Routing Problem

**Mohammadreza Nazari,** [1]   **Afshin Oroojlooy,** [1]   **Lawrence V. Snyder,** [1]   **Martin Takáč** [1]

## Abstract

We present an end-to-end framework for solving Vehicle Routing Problem (VRP) using deep reinforcement learning. In this approach, we train a single model that finds near-optimal solutions for problem instances sampled from a given distribution, only by observing the reward signals and following feasibility rules. Our model represents a parameterized stochastic policy, and by applying a policy gradient algorithm to optimize its parameters, the trained model produces the solution as a sequence of consecutive actions in real time, without the need to re-train for every new problem instance. Our method is faster in both training and inference than a recent method that solves the Traveling Salesman Problem (TSP), with nearly identical solution quality. On the more general VRP, our approach outperforms classical heuristics on medium-sized instances in both solution quality and computation time (after training). Our proposed framework can be applied to variants of the VRP such as the stochastic VRP, and has the potential to be applied more generally to combinatorial optimization problems.

## 1. Introduction

The *Vehicle Routing Problem* (VRP) is a combinatorial optimization problems that has been studied in applied mathematics and computer science for decades and for which many exact and heuristic algorithms have been proposed. It is known to be considerably more computationally difficult than the *Traveling Salesman Problem* (TSP), which is a special case. In the simplest form of the VRP, a single capacitated vehicle is responsible for delivering items to multiple customer nodes; the vehicle must return to the depot to pick up additional items when it runs out. The objective is to optimize a set of routes, all beginning and

[1]Department of Industrial Engineering, Lehigh University, PA, USA. Correspondence to: Mohammadreza Nazari <mon314@lehigh.edu>.

ending at a given node, called the *depot*, in order to attain the maximum possible reward, which is often the negative of the total vehicle traversed distance or average service time. For an overview of the VRP, see Laporte (1992); Laporte et al. (2000); Golden et al. (2008); Toth & Vigo (2002).

The prospect of new algorithm discovery, without any hand-engineered reasoning, makes neural networks and reinforcement learning a compelling choice that has the potential to be an important milestone on the path of approaching these problems. In this work, we develop a framework for solving a wide variety of combinatorial optimization problems using *Deep Reinforcement Learning* (DRL) and show how it can be applied to solve the VRP. For this purpose, we consider the Markov Decision Process (MDP) formulation of the problem, in which the optimal solution can be viewed as a sequence of decisions. This allows us to use DRL to produce near-optimal solutions by increasing the probability of decoding "desirable" sequences. A naive approach is to find a problem-specific solution by considering every instance separately. Obviously, this approach is not competitive with other algorithms in terms of either quality of solutions or runtime since there should be many trajectories sampled from one MDP to produce a near-optimal solution. Moreover, the learned policy does not apply to instances other than the one that was used in the training; with a small perturbation of the problem setting, we need to rebuild the policy from scratch.

Rather than focusing on training a separate model for every problem instance, we propose a structure that performs well on any problem sampled from a given distribution. One can view the trained model as a black-box heuristic (or a meta-algorithm) which generates a high-quality solution in a reasonable amount of time. Once the trained model is available, it can be used many times, without needing to re-train for the new problems as long as they are generated from the training distribution.

This work is originally motivated by the recent work by Bello et al. (2016). We have generalized their framework to include a wider range of combinatorial optimization problems such as the VRP. Bello et al. (2016) propose the use of a Pointer Network (Vinyals et al., 2015) to decode the solution. One major issue that prohibits the direct use of their approach for the VRP is that it assumes the system

is static over time. In contrast, in the VRP, the demands change over time in the sense that once a node has been visited its demand becomes, effectively, zero. To overcome this, we propose an alternate approach—which is actually simpler than the Pointer Network approach—that can efficiently handle both the static and dynamic elements of the system. Our model consists of a recurrent neural network (RNN) decoder coupled with an attention mechanism. At each time step, the embeddings of the static elements are the input to the RNN decoder, and the output of the RNN and the dynamic element embeddings are fed into an attention mechanism, which forms a distribution over the feasible inputs that can be chosen at the next decision point. When applied to the TSP, a special case of the VRP in which there is only a single route to optimize, our method produces solutions of similar quality as those reported by Bello et al. (2016), but is significantly faster in both training and inference time.

The proposed framework is appealing since we utilize a self-driven learning procedure that only requires the reward calculation based on the generated outputs; as long as we can observe the reward and verify the feasibility of a generated sequence, we can learn the desired meta-algorithm. For instance, if one does not know how to solve the VRP but can compute the cost of a given solution, then one can provide the signal required for solving the problem using our method. Unlike most classical heuristic methods, it is robust to problem changes, meaning that when the inputs change in any way, it can automatically adapt the solution. Using classical heuristics for VRP, the entire distance matrix must be recalculated and the system must be re-optimized from scratch, which is often impractical, especially if the problem size is large. In contrast, our proposed framework does not require an explicit distance matrix, and only one feed-forward pass of the network will update the routes based on the new data.

## 2. Background

Before presenting our framework, we briefly review some background that is closely related to our work.

### 2.1. Sequence-to-Sequence Models

*Sequence-to-Sequence* models (Sutskever et al., 2014; Vinyals et al., 2015; Luong et al., 2015) are useful in tasks for which a mapping from one sequence to another is required. They have been extensively studied in the field of neural machine translation over the past several years, and there are numerous variants of these models. The general architecture, which is almost the same among different versions, consists of two RNN networks, called the encoder and decoder. An encoder network reads through the input sequence and stores the knowledge in a fixed-size vector representation (or a sequence of vectors); then, a decoder converts the encoded information back to an output sequence.

In the vanilla Sequence-to-Sequence architecture (Sutskever et al., 2014), the source sequence appears only once in the encoder and the entire output sequence is generated based on one vector (i.e., the last hidden state of the encoder RNN). Other extensions, for example Bahdanau et al. (2015), illustrate that the source information can be used more wisely to increase the amount of information during the decoding steps. In addition to the encoder and decoder network, they employ another neural network, namely an *attention mechanism* that *attends* to the entire encoder RNN states. This mechanism allows the decoder to focus on the important locations of the source sequence and use the relevant information during decoding steps for producing "better" output sequences. Recently, the concept of attention has been a popular research idea due to its capability to align different objects, e.g., in computer vision (Chen et al., 2015; Xiao et al., 2015; Xu et al., 2015; Hong et al., 2016) and neural machine translation (Bahdanau et al., 2015; Jean et al., 2015; Luong et al., 2015). In this study, we also employ a special attention structure for policy representation. See Section 3.3 for a detailed discussion of the attention mechanism.

### 2.2. Neural Combinatorial Optimization

Over the last several years, multiple methods have been developed to tackle combinatorial optimization problems by using recent advances in artificial intelligence. The first attempt was proposed by Vinyals et al. (2015), who introduce the concept of a *Pointer Network*, a model originally inspired by sequence-to-sequence models. Because it is invariant to the length of the encoder sequence, the Pointer Network enables the model to apply to combinatorial optimization problems, where the output sequence length is determined by the source sequence. They use the Pointer Network architecture in a supervised fashion to find near-optimal TSP tours from ground truth optimal (or heuristic) solutions. This dependence on supervision prohibits the Pointer Network from finding better solutions than the ones provided during the training.

Closest to our approach, Bello et al. (2016) address this issue by developing a neural combinatorial optimization framework that uses RL to optimize a policy modeled by a Pointer Network. Using several classical combinatorial optimization problems such as TSP and the knapsack problem, they show the effectiveness and generality of their architecture.

On a related topic, Dai et al. (2017) solve optimization problems over graphs with using a graph embedding structure (Dai et al., 2016) and a DQN algorithm (Mnih et al., 2015). Even though VRP can be represented by a graph with weighted nodes and edges, their proposed model does not directly apply since in VRP, a particular node might be

visited multiple times.

Next, we introduce our model, which is a simplified version of the Pointer Network.

## 3. The Model

In this section, we formally define the problem and our proposed framework. Let us consider a generic combinatorial optimization problem with a given set of inputs $X \doteq \{x^i, i = 1, \cdots, M\}$. We allow some of the elements of each input to change between the decoding steps, which is, in fact, the case in many problems such as the VRP. The dynamic elements might be an artifact of the decoding procedure itself, or they can be imposed by the environment. For example, in the VRP, the remaining customer demands change over time as the vehicle visits the customer nodes; or we might consider a variant in which new customers arrive or adjust their demand values over time, independent of the vehicle decisions. Formally, we represent each input $x^i$ by a sequence of tuples $\{x_t^i = (s^i, d_t^i), t = 0, \cdots, T\}$, where $s^i$ and $d_t^i$ are the static and dynamic elements of the input, respectively, and can themselves be tuples. For instance, in the VRP, $s^i$ corresponds to the 2-dimensional coordinate of customer $i$'s location and $d_t^i$ is its demand at time $t$. We will denote the set of all input states at a fixed time $t$ with $X_t$.

We start from an arbitrary input $y_0 \in X_0$. At every decoding time $t = 0, 1, \cdots$, we choose $y_{t+1}$ from the available inputs $X_t$ and continue until a termination condition is satisfied. The termination condition is problem-specific, showing that the generated sequence satisfies the feasibility constraints. For instance, the termination criterion in the TSP is that all cities are visited; or in the VRP that we consider in this work, the terminating condition is that there is no more demand to satisfy. This process will generate a sequence of length $T'$, $Y = \{y_t, t = 0, ..., T'\}$, possibly with a different sequence length compared to the inputs and different dynamic elements. We also use the notation $Y_t$ to denote the decoded sequence up to time $t$, i.e. $Y_t = \{y_0, \cdots, y_t\}$. We are interested in finding a stochastic policy $\pi$ which generates the sequence $Y$ in a way that minimizes a loss objective while satisfying the problem constraints. The optimal policy $\pi^*$ will generate the optimal solution with probability 1. Our goal is to make $\pi$ as close to $\pi^*$ as possible. Similar to Sutskever et al. (2014), we use the probability chain rule to decompose the probability of generating sequence $Y$, i.e., $P(Y|X_0)$, as follows:

$$P(Y|X_0) = \prod_{t=0}^{T'} P(y_{t+1}|Y_t, X_t), \qquad (1)$$

and

$$X_{t+1} = f(y_{t+1}, X_t) \qquad (2)$$

is a recursive update of the problem representation with the state transition function $f$. Each component in right-hand side of (1) is computed by the attention mechanism, i.e.,

$$P(y_{t+1}|Y_t, X_t) = \text{softmax}(g(h_t, X_t)), \qquad (3)$$

where $g$ is an affine function that outputs an input-sized vector, and $h_t$ is the state of the RNN decoder that summarizes the information of previously decoded steps $y_0, \cdots, y_t$. We will describe the details of our proposed attention mechanism in Section 3.3.

**Remark**: This model can handle combinatorial optimization problems in both a more classical static setting as well as in dynamically changing ones. In static combinatorial optimization, $X_0$ fully defines the problem that we are trying to solve. For example, in the VRP, $X_0$ includes all customer locations as well as their demands, and the depot location; then, the remaining demands are updated with respect to the vehicle destination and its load. With this consideration, often there exists a well-defined Markovian transition function $f$, as defined in (2), which is sufficient to update the dynamics between decision points. However, our model can also be applied to problems in which the state transition function is unknown and/or is subject to external noise, since the training does not explicitly make use of the transition function. This transition function helps simulate the environment that the training algorithm interacts with. See Section 5.4 for an example of how to apply the model to the stochastic VRP.

### 3.1. Limitations of Pointer Networks

Although the framework proposed by Bello et al. (2016) works well on problems such as the knapsack problem and TSP, it is not applicable to more complicated combinatorial optimization problems in which the system representation varies over time, such as VRP. Bello et al. (2016) feed a random sequence of inputs to the RNN encoder. Figure 1 illustrates with an example why using the RNN in the encoder is restrictive. Suppose that at the first decision step, the policy sends the vehicle to customer 1, and as a result, its demand is satisfied, i.e., $d_0^1 \neq d_1^1$. Then in the second decision step, we need to re-calculate the whole network with the new $d_1^1$ information in order to choose the next customer. The dynamic elements complicate the forward pass of the network since there should be encoder/decoder updates when an input changes. The situation is even worse during back-propagation to accumulate the gradients since we need to remember when the dynamic elements changed. In order to resolve this complication, we require the model to be *invariant to the input sequence* so that changing the order of any two inputs does not affect the network. In Section 3.2, we present a simple network that satisfies this property.
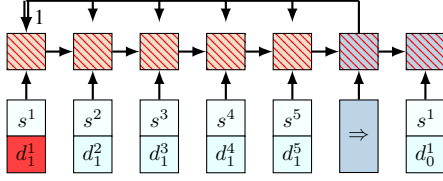
Figure 1. Limitation of the Pointer Network. After a change in dynamic elements ($d_1^1$ in this example), the whole Pointer Network should be updated for computing the probabilities in the next decision point.

### 3.2. The Proposed Neural Network Model

We argue that RNN encoder adds extra complication to the encoder but is actually not necessary, and the approach can be made much more general by omitting it. RNNs are necessary only when the inputs convey sequential information; e.g., in text translation the combination of words and their relative position must be captured in order for the translation to be accurate. But the question here is, *why do we need to have them in the encoder for combinatorial optimization problems when there is no meaningful order in the input set?* As an example, in the TSP, the inputs are the set of unordered city locations, and their order is not meaningful; any random permutation contains the same information as the original inputs. Therefore, in our model, we simply leave out the encoder RNN and directly use the embedded inputs instead of the RNN hidden states. By this modification, many of the computational complications disappear, without decreasing the model's efficiency.

As illustrated in Figure 2, our model is composed of two main components: *(i)* a set of embeddings that maps the inputs into a $D$-dimensional vector space. We might have multiple embeddings corresponding to different elements of the input, but they are shared among the inputs. We use 1-dimensional convolution layers for the embedding, in which the in-width is the input length, the number of filters is $D$, and the number of in-channels is the number of elements of $x$; and *(ii)* a decoder that points to an input at every decoding step. As is common in the literature (Bahdanau et al., 2015; Sutskever et al., 2014; Cho et al., 2014), we use RNN to model the decoder network. For similar reasons as in the discussion in Section 3.1, notice that only the static elements are the inputs to the decoder network; we only use the dynamic elements in the attention layer, described next.

### 3.3. Attention Mechanism

An attention mechanism is a differentiable structure for addressing different parts of the input. Figure 3 illustrates the attention mechanism employed in our method. At decoder step $i$, we utilize a context-based attention mechanism, similar to Vinyals et al. (2015), which extracts the relevant
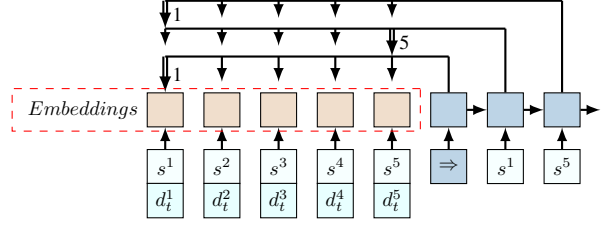


Figure 2. Our proposed model. The embedding layer, on the left, maps the inputs to a high-dimensional vector space. On the right, an RNN decoder stores the information of the decoded sequence. Then, the RNN hidden state and embedded input produce a probability distribution over the next input using the attention mechanism.

information from the inputs using a variable-length alignment vector $a_t$. In words, $a_t$ specifies how much every input data point might be relevant in the next decoding step $t$. Recall that for simplicity of calculations, the dynamic elements are only used in the attention layer.
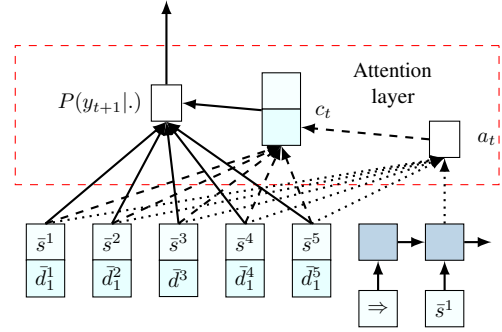


Figure 3. The proposed attention mechanism.

Let $\bar{x}_t^i = (\bar{s}^i, \bar{d}_t^i)$ be the embedded input $i$, and $h_t \in \mathbb{R}^D$ be the memory state of the RNN cell at decoding step $t$. The alignment vector $a_t$ is then computed as

$$a_t = a_t(\bar{x}_t^i, h_t) = \text{softmax}(u_t), \quad (4)$$

where

$$u_t^i = v_a^T \tanh\left(W_a[\bar{x}_t^i; h_t]\right). \quad (5)$$

Here ";" means the concatenation of two vectors. We compute the conditional probabilities by combining the context vector $c_t$, computed as

$$c_t = \sum_{m=1}^{M} a_t^i \bar{x}_t^i, \quad (6)$$

with the embedded inputs, and then normalizing the values with the softmax function, as follows:

$$P(y_{t+1}|Y_t, X_t) = \text{softmax}(\tilde{u}_t^i), \quad (7)$$

$$\tilde{u}_t^i = v_c^T \tanh\left(W_c[\bar{x}_t^i; c_t]\right). \quad (8)$$

In (4)–(8), $v_a$, $v_c$, $W_a$ and $W_c$ are trainable variables.

**Remark**: *Model Symmetry*: Vinyals et al. (2016) discuss an extension of sequence-to-sequence models where they empirically demonstrate that in tasks with no obvious input sequence, such as sorting, the order in which the inputs are fed into the network matter. A similar concern arises when using Pointer Networks for combinatorial optimization problems. However, the model proposed in this paper does not suffer from such a complication since the embeddings and the attention mechanism are invariant to the input order.

## 4. Training method

To train the network, we use well-known policy gradient approaches. To use these methods, we parameterize the stochastic policy $\pi$ with parameters $\theta$. Policy gradient methods use an estimate of the gradient of the expected return with respect to the policy parameters to iteratively improve the policy. We utilize the REINFORCE method, similar to Bello et al. (2016) for solving the TSP and VRP, and A3C (Mnih et al., 2016) for the stochastic VRP. In both of these algorithms, there are two networks: (*i*) an actor network that predicts a probability distribution over the next action at any given decision step, and (*ii*) a critic network that estimates the reward for any problem instance from a given state. In this paper, the critic network, first, uses the output probabilities of the actor network to compute a weighted sum of the embedded inputs, and then, it has two hidden layers: one dense layer with ReLU activation and another linear one with single output. In this section, we elaborate the details of the REINFORCE algorithm; see the details of A3C in Appendix D.

Let us consider a family of problems, denoted by $\mathcal{M}$, and a probability distribution over them, denoted by $\Phi_{\mathcal{M}}$. During the training, the problem instances are generated according to distribution $\Phi_{\mathcal{M}}$. We also use the same distribution in the inference to produce test examples. Algorithm 1 summarizes the training algorithm. We have two neural networks with weight vectors $\theta$ and $\phi$ associated with the actor and critic, respectively. We draw $N$ sample problems from $\mathcal{M}$ and use Monte Carlo simulation to produce feasible sequences with respect to the current policy $\pi_\theta$. We adopt the superscript $n$ to refer to the variables of the $n$th instance. After termination of the decoding in all $N$ problems, we compute the corresponding rewards as well as the policy gradient in step 14 to update the actor network. In this step, $V(X_0^n; \phi)$ is the the reward approximation for instance problem $n$ that will be calculated from the critic network. We also update the critic network in step 15 in the direction of reducing the difference between the expected rewards with the observed ones during Monte Carlo roll-outs.

---

**Algorithm 1** REINFORCE Algorithm

1: initialize the actor network with random weights $\theta$ and critic network with random weights $\phi$
2: **for** $iteration = 1, 2, \cdots$ **do**
3:    reset gradients: $d\theta \leftarrow 0$, $d\phi \leftarrow 0$
4:    sample $N$ instances according to $\Phi_{\mathcal{M}}$
5:    **for** $n = 1, \cdots, N$ **do**
6:       initialize step counter $t \leftarrow 0$
7:       **repeat**
8:          choose $y_{t+1}^n$ according to the distribution $P(y_{t+1}^n | Y_t^n, X_t^n)$
9:          observe new state $X_{t+1}^n$
10:          $t \leftarrow t + 1$
11:       **until** termination condition is satisfied
12:       compute reward $R^n = R(Y^n, X_0^n)$
13:    **end for**
14:    $d\theta \leftarrow \frac{1}{N} \sum_{n=1}^N (R^n - V(X_0^n; \phi)) \nabla_\theta \log P(Y^n | X_0^n)$
15:    $d\phi \leftarrow \frac{1}{N} \sum_{n=1}^N \nabla_\phi (R^n - V(X_0^n; \phi))^2$
16:    update $\theta$ using $d\theta$ and $\phi$ using $d\phi$.
17: **end for**

---

## 5. Experimental Results

Using TSP as the test-bed, the first experiment of this section is designed to validate the performance of the proposed method compared to the model of Bello et al. (2016). In the second experiment, we study the capacitated VRP and compare its results with two existing heuristic algorithms. Finally, we address the extensions to other VRP variants with demonstrating an experiment on stochastic VRP.

### 5.1. TSP Results

To evaluate the effectiveness of our framework, we compare the route lengths of the TSP solutions obtained by our framework with those given by model of Bello et al. (2016) for random instances with 20, 50, and 100 nodes. In the training phase, we generate 1,000,000 TSP instances for each problem size, and use them in training for 20 epochs. The city locations are chosen uniformly from the unit square $[0, 1] \times [0, 1]$. We use the same data distribution to produce instances for the testing phase.

Table 1 summarizes the result for different TSP sizes using the greedy decoder in which at every decoding step, the city with the highest probability is chosen as the destination. The results are averaged over 1000 instances. The first column is the average TSP tour length using our proposed architecture, the second column is the result of our implementation of Bello et al. (2016), and the optimal tour lengths are reported in the last column. A comparison of the first two columns suggests that there is almost no difference between the performance of two approaches. In fact, the RNN encoder

of the Pointer Network learns to convey no information to the next steps, i.e., $h_t = f(x_t)$. On the other hand, our approach is around 60% faster in both training and inference, since it has two fewer RNNs—one in the encoder of actor network and another in the encoder of critic network. Table 2 summarizes the training times for one epoch of the training and the time-savings that we gain by eliminating the encoder RNN.

*Table 1.* Average tour length for TSP.

| Task | Our Framework (Greedy) | Pointer-RL (Greedy) | Optimal |
|------|------------------------|---------------------|---------|
| TSP20 | 3.97 | 3.96 | 3.82 |
| TSP50 | 6.08 | 6.05 | 5.70 |
| TSP100 | 8.44 | 8.45 | 7.77 |

*Table 2.* TSP training time for one epoch (in minutes).

| Task | Our Framework (Greedy) | Pointer-RL (Greedy) | % Time Saving |
|------|------------------------|---------------------|---------------|
| TSP20 | 22.18 | 50.33 | 55.9% |
| TSP50 | 54.10 | 147.25 | 63.3% |
| TSP100 | 122.10 | 300.73 | 59.4% |

## 5.2. VRP Experimental Design

Many variants of the VRP have been extensively studied in operations research (See, for example, the reviews by Laporte (1992); Laporte et al. (2000), or the book by Toth & Vigo (2002) for different variants of the problem). In this section, we consider a specific capacitated version of the problem in which one vehicle with a limited capacity is responsible for delivering items to many geographically distributed customers with finite demands. When the vehicle's load runs out, it returns to the depot to refill. We will denote the vehicle's remaining load at time $t$ as $l_t$. The objective is to minimize the total route length while satisfying all of the customer demands.

We assume that the node locations and demands are randomly generated from a fixed distribution. Specifically, the customers and depot locations are randomly generated in the unit square $[0, 1] \times [0, 1]$. For simplicity, we assume that the demand of each node is a discrete number in $\{1, .., 9\}$, chosen uniformly at random. We note, however, that the demand values can be generated from any distribution, including continuous ones.

We assume that the vehicle is located at the depot at time 0, so the first input to the decoder is an embedding of the depot location. At each decoding step, the vehicle chooses from among the customer nodes or the depot to visit in the next step. After visiting customer node $i$, the demands and

vehicle load are updated as follows:

$$d_{t+1}^i = \max(0, d_t^i - l_t), \qquad (9)$$
$$d_{t+1}^k = d_t^k, \quad \text{for } k \neq i, \qquad (10)$$
$$l_{t+1} = \max(0, l_t - d_t^i), \qquad (11)$$

which is an explicit definition of the state transition function (2) for the VRP.

In this experiment, we have employed two different decoders: *(i)* greedy, in which at every decoding step, the node (either customer or depot) with the highest probability is selected as the next destination, and *(ii)* beam search (BS), which keeps track of the most probable paths and then chooses the one with the highest reward. Our results indicate that by applying the beam search algorithm, the quality of the solutions can be improved with only a slight increase in computation time. For faster training and generating feasible solutions, we have used a masking scheme which sets the log-probabilities of infeasible solutions to $-\infty$ or forces a solution if a particular condition is satisfied. In the VRP, we use the following masking procedures: *(i)* nodes with zero demand are not allowed to be visited; and *(ii)* all customer nodes will be masked if the vehicle's remaining load is exactly 0.

We use one layer of LSTM RNN in the decoder with a state size of 128. Each customer location is also embedded into a vector of size 128, shared among the inputs. We employ similar embeddings for the dynamic elements; the demand $d_t^i$ and the remaining vehicle load after visiting node $i$, $l_t - d_t^i$, are mapped to a vector in a 128-dimensional vector space and used in the attention layer. The variables in both actor and critic network are initialized with Xavier initialization (Glorot & Bengio, 2010). For training both networks, we use the Adam optimizer (Kingma & Ba, 2015) with learning rate $10^{-4}$. The batch size $N$ is 128, and we clip the gradients when their norm is greater than 2. We use dropout with probability 0.1 in the decoder LSTM. Moreover, we tried the entropy regularizer (Williams & Peng, 1991; Mnih et al., 2016), which has been shown to be useful in preventing the algorithm from getting stuck in local optima, but it does not show any improvement in our experiments; therefore, we do not use it in the results reported below.

On a single GPU K80, every 100 training steps of the VRP with 20 customer nodes takes approximately 35 seconds. Training for 20 epochs requires about 13.5 hours. The TensorFlow implementation of our code for both TSP and VRP will be publicly available.

## 5.3. VRP Results

In Tables 3 and 4, we compare the solutions found by using our framework with those obtained from the *Clarke-Wright*

*Table 3.* Average total tour length over a test set of size 1000. Standard deviations are in parentheses.

| Task | Vehicle Capacity | Our Framework | | | Clarke-Wright | | | Sweep | | |
|------|------|------|------|------|------|------|------|------|------|------|
| | | Greedy | BS(3) | BS(10) | Greedy | Rnd(5,5) | Rnd(10,10) | Basic | Rnd(5) | Rnd(10) |
| VRP10 | 20 | 4.80 | 4.72 | **4.65** | 5.06 | 4.86 | 4.80 | 5.42 | 5.07 | 5.00 |
| | | (0.83) | (0.80) | (0.79) | (0.85) | (0.82) | (0.82) | (0.95) | (0.87) | (0.87) |
| VRP20 | 30 | 6.51 | 6.45 | **6.34** | 7.22 | 6.89 | 6.81 | 7.59 | 7.17 | 7.08 |
| | | (0.84) | (0.85) | (0.80) | (0.90) | (0.84) | (0.82) | (0.93) | (0.85) | (0.84) |
| VRP50 | 40 | 11.32 | 11.21 | **11.08** | 12.85 | 12.35 | 12.25 | 13.61 | 13.09 | 12.96 |
| | | (1.27) | (1.30) | (1.27) | (1.33) | (1.27) | (1.25) | (1.23) | (1.12) | (1.12) |
| VRP100 | 50 | 17.12 | 17.00 | **16.86** | 19.72 | 19.09 | 18.96 | 21.01 | 20.47 | 20.33 |
| | | (1.90) | (1.88) | (1.87) | (1.92) | (1.85) | (1.85) | (1.51) | (1.41) | (1.39) |

*Table 4.* Solution times (in seconds).

| Task | Vehicle Capacity | Our Framework | | | Clarke-Wright | | | Sweep | | |
|------|------|------|------|------|------|------|------|------|------|------|
| | | Greedy | BS(3) | BS(10) | Greedy | Rnd(5,5) | Rnd(10,10) | Basic | Rnd(5) | Rnd(10) |
| VRP10 | 20 | 0.059 | 0.073 | 0.074 | 0.002 | 0.016 | 0.079 | 0.001 | 0.004 | 0.008 |
| VRP20 | 30 | 0.107 | 0.147 | 0.155 | 0.011 | 0.053 | 0.256 | 0.006 | 0.029 | 0.062 |
| VRP50 | 40 | 0.176 | 0.244 | 0.250 | 0.052 | 0.217 | 0.903 | 0.096 | 0.472 | 0.988 |
| VRP100 | 50 | 0.290 | 0.413 | 0.477 | 0.186 | 0.735 | 3.171 | 1.341 | 6.32 | 12.443 |

*savings* and *Sweep* heuristic algorithms (See Appendix C for more details of both heuristics). We run our test on multiple problem sizes with different vehicle capacities; for example, VRP10 consists of 10 customer nodes and one depot. The results are based on 1000 instances, sampled for each problem size.

Table 3 shows the average total lengths of the routes generated by our framework, using greedy and BS decoders, with the number inside the parentheses in the column header indicating the beam-width parameter. In addition, we also implemented a randomized version of both heuristic algorithms to improve the solution quality; for Clarke-Wright, the numbers inside the parentheses are the randomization depth and randomization iterations parameters ; and for Sweep, it is the number of random initial angles for grouping the nodes. We observe that the average total length of the solutions found by our method using various decoders outperforms the heuristic algorithms. We also see that using the beam search decoder significantly improves the solution while only adding a small computational cost in run-time. This table also provides the standard deviation in parenthesis (below each average). These numbers are close to each other between our framework and Clarke-Wright Savings algorithms, suggesting that the difference in the average tour length is not due to pathological examples in which one of the algorithms fails. Table 4 presents the solution time comparison. Even though the greedy Clark-Wright and basic Sweep heuristics are the fastest for small instances, they do not provide competitive solutions. Moreover, for larger problems, our framework is faster than the randomized heuristics.

We find that training without an embedding layer always yields an inferior solution. One possible explanation is that the policy is able to extract useful features from the high-dimensional input representations much more efficiently. Recall that our embedding is an affine transformation, so it does not necessarily keep the embedded input distances proportional to the original 2-dimensional Euclidean distances.

Figure 4 shows the log of the ratio of solution times to the number of customer nodes. We observe that this ratio stays almost the same for our framework with different decoders. In contrast, the log of the run time for the Clarke-Wright and Sweep heuristics increases faster than linearly with the number of nodes. This observation is one motivation for applying our framework to more general combinatorial problems, since it suggests that our method scales well.

### 5.4. Extension to Other VRPs: Stochastic VRP

The proposed framework can be extended easily to problems with multiple depots; one only needs to construct the corresponding state transition function and masking procedure. It is also possible to incorporate various side constraints: soft constraints can be applied by penalizing the rewards, or hard constraints such as time windows can be enforced through a masking scheme. However, designing such a scheme might be a challenging task, possibly harder than solving the optimization problem itself.

This framework can also be extended to real-time services including on-demand deliveries and taxis. A major difficulty of planning in these systems is that the schedules are not defined in beforehand, and one needs to deal with
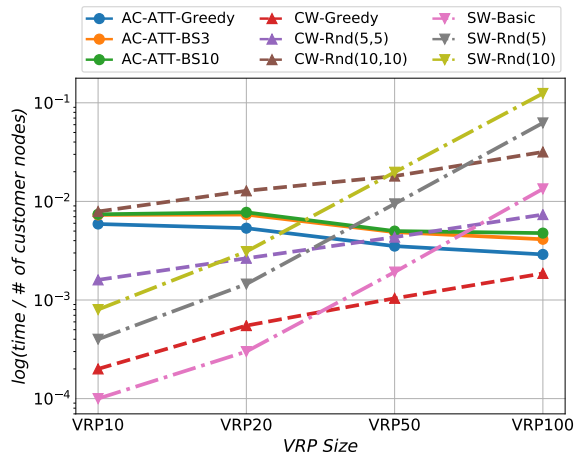
*Figure 4.* Log of ratio of solution time to the number of customer nodes using different algorithms.

various demand realization on the fly. Next, we design a simulated experiment to illustrate the performance of the framework on the *stochastic VRP*. We consider an instance of the stochastic VRP in which customers with random demands arrive at the system according to a Poisson process; without loss of generality we assume the process has rate 1. Similar to previous experiments, we choose each new customer's location uniformly on the unit square and its demand to a discrete number in $\{1, \cdots, 9\}$. We fix the depot position to $[0.5, 0.5]$. A vehicle is required to satisfy as much demand as possible in a time horizon with length 100 time units. To make the system stable, we assume that each customer cancels its demand after having gone unanswered for 5 time units. The vehicle moves with speed 0.1 per time unit. Obviously, this is a continuous-time system, but we view it as a discrete-time MDP where the vehicle can make decisions at either the times of customer arrivals or after the time when the vehicle reaches a node.

The network and its hyper-parameters in this experiment are the same as in the previous experiments. One major difference is the RL training method, where we use asynchronous advantage actor-critic (A3C) (Mnih et al., 2016) with one-step reward accumulation. The details of the training method are described in Appendix D. The other difference is that instead of using masking, at every time step, the input to the network is a set of available locations which consists of the customers with positive demand, the depot, and the vehicle's current location; the latter decision allows the vehicle to stop at its current position, if necessary. We also add the *time-in-system* of customers as a dynamic element to the attention mechanism; it will allow the training process to learn customer abandonment behavior.

We compare our results with three other strategies: *(i) Random*, in which the next destination is randomly selected

from the available nodes; *(ii) Largest-Demand*, in which the customer with maximum demand will be chosen as the next destination; and *(iii) Max-Reachable*, in which the vehicle chooses the node with the highest demand while making sure that the demand will remain valid until the vehicle reaches the node. In all strategies, we force the vehicle to route to the depot and refill when its load is zero.

Table 5 summarizes the average demand satisfied, and the percentage of the total demand that this represents, under the various strategies, averaged over 100 test instances. We observe that A3C outperforms the other strategies. Even though A3C does not know any information about the problem structure, it is able to perform better than the Max-Reachable strategy, which uses customer abandonment information.

*Table 5.* Satisfied demand under different strategies.

| Method | Random | Largest-Demand | Max-Reachable | A3C |
|---|---|---|---|---|
| Avg. Dem. | 24.83 | 75.11 | 88.60 | 112.21 |
| % satisfied | 5.4% | 16.6% | 19.6% | 28.8% |

## 6. Discussion and Conclusion

We expect that the proposed architecture has a large potential to be used in real-world problems with further improvements. One possible approach that as not been studied here is to replace the REINFORCE or A3C algorithm with other policy-based or value-based RL algorithms. Also noting that the proposed algorithm is not limited to TSP or VRP, it will be an important future research to apply it to other combinatorial optimization problems such as bin-packing, job-shop, flow-shop.

This method is quite appealing because if we generate a new VRP instance with the same number of nodes and vehicle capacity, and the same location and demand distributions as the ones that we used during training, then the trained policy will work well, and we can solve the problem right away, without retraining for every new instance. As long as we approximate the generating distribution of the problem, the framework can be applied. The only requirement is a verifier to find feasible solutions and also a reward signal to demonstrate how well the policy is working.

Unlike many classical heuristics, our proposed method scales well with increasing problem size, and has a superior performance with competitive solution-time. It doesn't require a distance matrix calculation which might be computationally cumbersome, especially in dynamically changing VRPs. We also illustrate the performance of the algorithm on a much more complicated stochastic version of the VRP.

## Acknowledgment

## References

Bahdanau, Dzmitry, Cho, Kyunghyun, and Bengio, Yoshua. Neural machine translation by jointly learning to align and translate. *In International Conference on Learning Representations*, 2015.

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

Chen, Kan, Wang, Jiang, Chen, Liang-Chieh, Gao, Haoyuan, Xu, Wei, and Nevatia, Ram. Abc-cnn: An attention based convolutional neural network for visual question answering. *arXiv preprint arXiv:1511.05960*, 2015.

Cho, Kyunghyun, Van Merriënboer, Bart, Gulcehre, Caglar, Bahdanau, Dzmitry, Bougares, Fethi, Schwenk, Holger, and Bengio, Yoshua. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *Conference on Empirical Methods in Natural Language Processing*, 2014.

Clarke, Geoff and Wright, John W. Scheduling of vehicles from a central depot to a number of delivery points. *Operations research*, 12(4):568–581, 1964.

Dai, Hanjun, Dai, Bo, and Song, Le. Discriminative embeddings of latent variable models for structured data. In *International Conference on Machine Learning*, pp. 2702–2711, 2016.

Dai, Hanjun, Khalil, Elias B, Zhang, Yuyu, Dilkina, Bistra, and Song, Le. Learning combinatorial optimization algorithms over graphs. *Advances in Neural Information Processing Systems*, 2017.

Glorot, Xavier and Bengio, Yoshua. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 249–256, 2010.

Golden, Bruce L, Raghavan, Subramanian, and Wasil, Edward A. *The Vehicle Routing Problem: Latest Advances and New Challenges*, volume 43. Springer Science & Business Media, 2008.

Hong, Seunghoon, Oh, Junhyuk, Lee, Honglak, and Han, Bohyung. Learning transferrable knowledge for semantic segmentation with deep convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3204–3212, 2016.

Jean, Sébastien, Cho, Kyunghyun, Memisevic, Roland, and Bengio, Yoshua. On using very large target vocabulary for neural machine translation. 2015.

Kingma, Diederik P and Ba, Jimmy. Adam: A method for stochastic optimization. In *International Conference on Machine Learning*, 2015.

Laporte, Gilbert. The vehicle routing problem: An overview of exact and approximate algorithms. *European journal of operational research*, 59(3):345–358, 1992.

Laporte, Gilbert, Gendreau, Michel, Potvin, Jean-Yves, and Semet, Frédéric. Classical and modern heuristics for the vehicle routing problem. *International transactions in operational research*, 7(4-5):285–300, 2000.

Luong, Minh-Thang, Pham, Hieu, and Manning, Christopher D. Effective approaches to attention-based neural machine translation. *Conference on Empirical Methods in Natural Language Processing*, 2015.

Mnih, Volodymyr, Kavukcuoglu, Koray, Silver, David, Rusu, Andrei A, Veness, Joel, Bellemare, Marc G, Graves, Alex, Riedmiller, Martin, Fidjeland, Andreas K, Ostrovski, Georg, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

Mnih, Volodymyr, Badia, Adria Puigdomenech, Mirza, Mehdi, Graves, Alex, Lillicrap, Timothy, Harley, Tim, Silver, David, and Kavukcuoglu, Koray. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*, pp. 1928–1937, 2016.

Snyder, Lawrence V and Shen, Zuo-Jun Max. *Fundamentals of Supply Chain Theory*. John Wiley & Sons, 2 edition, 2018.

Sutskever, Ilya, Vinyals, Oriol, and Le, Quoc V. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pp. 3104–3112, 2014.

Toth, Paolo and Vigo, Daniele. *The Vehicle Routing Problem*. SIAM, 2002.

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

Vinyals, Oriol, Bengio, Samy, and Kudlur, Manjunath. Order matters: Sequence to sequence for sets. 2016.

Williams, Ronald J and Peng, Jing. Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 3(3):241–268, 1991.

Wren, Anthony and Holliday, Alan. Computer scheduling of vehicles from one or more depots to a number of delivery points. *Operational Research Quarterly*, pp. 333–344, 1972.

Xiao, Tianjun, Xu, Yichong, Yang, Kuiyuan, Zhang, Jiaxing, Peng, Yuxin, and Zhang, Zheng. The application of two-level attention models in deep convolutional neural network for fine-grained image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 842–850, 2015.

Xu, Kelvin, Ba, Jimmy, Kiros, Ryan, Cho, Kyunghyun, Courville, Aaron, Salakhudinov, Ruslan, Zemel, Rich, and Bengio, Yoshua. Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning*, pp. 2048–2057, 2015.

## A. Sample VRP Solutions

Figure 5 illustrates sample VRP20 and VRP50 instances decoded by the trained model. The greedy and beam-search decoders were used to produce the figures in the top and bottom rows, respectively. It is evident that these solutions are not optimal. For example, in part (a), one of the routes crosses itself, which is never optimal in Euclidean VRP instances. Another similar suboptimality is evident in part (c) to make the total distance shorter. However, the figures illustrate how well the policy model has understood the problem structure. It tries to satisfy demands at nearby customer nodes until the vehicle load is small. Then, it automatically comprehends that visiting further nodes is not the best decision, so it returns to the depot and starts a new tour. One interesting behavior that the algorithm has learned can be seen in part (c), in which the solution reduces the cost by making a partial delivery; in this example, we observe that the red and blue tours share a customer node with demand 8, each satisfying a portion of its demand; in this way, we are able to meet all demands without needing to initiate a new tour. We also observe how using the beam-search decoder produces further improvements; for example, as seen in parts (b)–(c), it reduces the number of times when a tour crosses itself; or it reduces the number of tours required to satisfy all demands as is illustrated in (b).

## B. Attention Mechanism Visualization

In order to illustrate how the attention mechanism is working, we relocated customer node 0 to different locations and observed how it affects the selected action. Figure 6 illustrates the attention in initial decoding step for a VRP10 instance drawn in part (a). Specifically, in this experiment, we let the coordinates of node 0 equal $\{0.1 \times (i, j), \forall i, j \in \{1, \cdots, 9\}\}$. In parts (b)-(d), the small bottom left square corresponds to the case where node 0 is located at [0.1,0.1] and the others have a similar interpretation. Each small square is associated with a color ranging from black to white, representing the probability of selecting the corresponding node at the initial decoding step. In part (b), we observe that if we relocate node 0 to the bottom-left of the plane, there is a positive probability of directly going to this node; otherwise, as seen in parts (c) and (d), either node 2 or 9 will be chosen with high probability. We do not display the probabilities of the other points since there is a near-0 probability of choosing them, irrespective of the location of node 0. A video demonstration of the model and attention mechanism is available online at https://streamable.com/gadhf.

## C. Heuristic Algorithms

In this section, we briefly describe the two heuristic algorithms that we used as benchmarks for the VRP. More details and examples of these algorithms can be found at Snyder & Shen (2018).

### C.1. Clarke-Wright Savings Heuristic

The Clarke-Wright savings heuristic (Clarke & Wright, 1964) is one of the best-known heuristics for the VRP. Let $\mathcal{N} \doteq \{1, \cdots, N\}$ be the set of customer nodes, and 0 be the depot. The distance between nodes $i$ and $j$ is denoted by $c_{ij}$, and $c_{0i}$ is the distance of customer $i$ from the depot. Algorithm 2 describes a randomized version of the heuristic. The basic idea behind this algorithm is that it initially considers a separate route for each customer node $i$, and then reduces the total cost by iteratively merging the routes. Merging two routes by adding the edge $(i, j)$ reduces the total distance by $s_{ij}$, so the algorithm prefers mergers with the highest savings. When $M = R = 1$, this algorithm is equivalent to the original Clarke-Wright savings heuristic, in which case, the feasible merger with the highest savings will be selected. By allowing $M, R > 1$, we introduce randomization, which can improve the performance of the algorithm further. In particular, Algorithm 2 chooses randomly from the $r \in \{1, \cdots, R\}$ best feasible mergers. Then, for each $r$, it solves the problem $m \in \{1, \cdots, M\}$ times, and returns the solution with the shortest total distance. In this paper, we refer to hyper-parameters $R$ and $M$ as the *randomization depth* and *randomization iteration*, respectively.

---

**Algorithm 2** Randomized Clarke-Wright Savings Heuristic

---

1: compute savings $s_{ij}$, where

$$s_{ij} = c_{i0} + c_{0j} - c_{ij} \qquad i, j \in \mathcal{N}, \, i \neq j$$
$$s_{ii} = 0 \qquad\qquad\qquad i \in \mathcal{N}$$

2: **for** $r = 1, \cdots, R$ **do**
3:     **for** $m = 1, \cdots, M$ **do**
4:         place each $i \in \mathcal{N}$ in its own route
5:         **repeat**
6:           find $k$ feasible mergers $(i, j)$ with the highest $s_{ij} > 0$, satisfying the following conditions:
              i) $i$ and $j$ are in different routes
              ii) both $i$ and $j$ are adjacent to the depot
              iii) combined demand of routes containing $i$ and $j$ is $\leq$ vehicle capacity
7:           choose a random $(i, j)$ from the feasible mergers, and combine the associated routes by replacing $(i, 0)$ and $(0, j)$ with $(i, j)$
8:         **until** no feasible merger is left
9:     **end for**
10: **end for**
11: **Return**: route with the shortest length

---

Path length = 5.72          Path length = 6.36          Path length = 10.72

Path length = 5.62          Path length = 6.06          Path length = 10.70

(a) Example 1 for VRP20; capacity 30      (b) Example 2 for VRP20; capacity 30      (c) VRP50; capacity 40
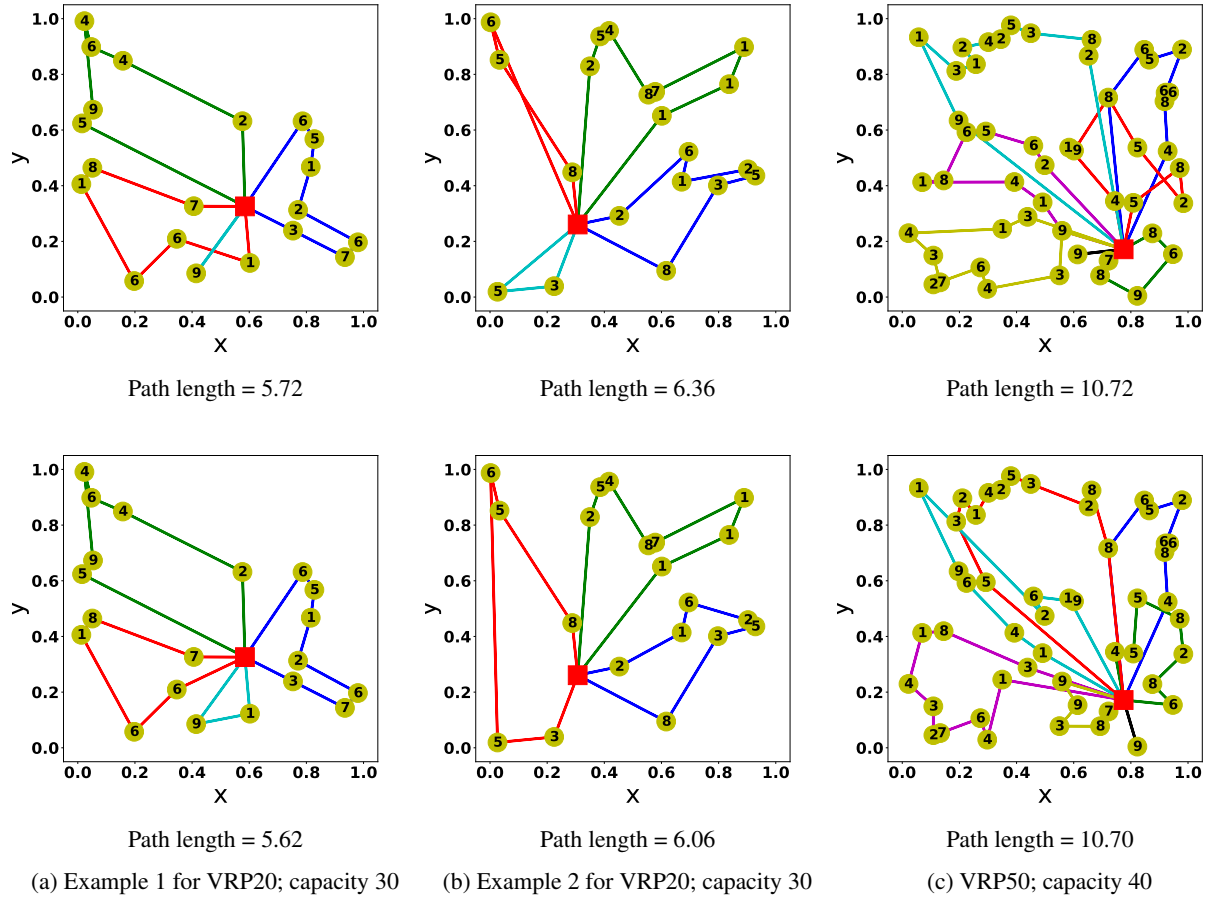
*Figure 5.* Sample decoded solutions for VRP20 and VRP50 using greedy (in top row) and beam-search (bottom row) decoder. The numbers inside the nodes are the demand values.
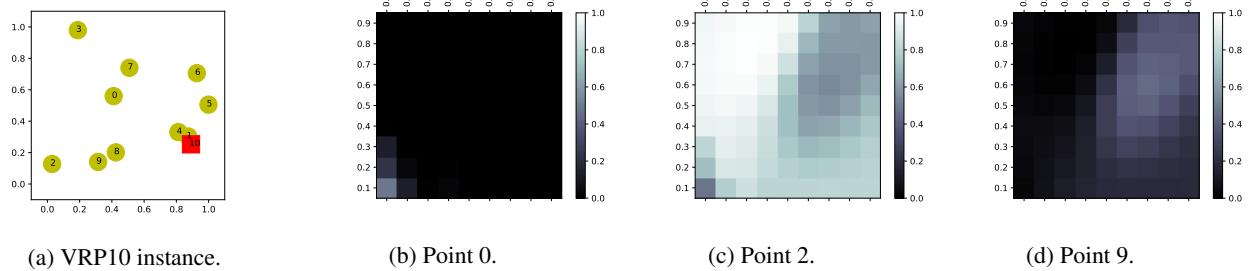


(a) VRP10 instance.          (b) Point 0.          (c) Point 2.          (d) Point 9.

*Figure 6.* Illustration of attention mechanism at decoding step 0. The problem instance is illustrated in part (a) where the nodes are labeled with a sequential number; labels 0-9 are the customer nodes and 10 is the depot. We place node 0 at different locations and observe how it affects the probability distribution of choosing the first action, as illustrated in parts (b)–(d).

## C.2. Sweep Heuristic

The sweep heuristic (Wren & Holliday, 1972) solves the VRP by breaking it into multiple TSPs. By rotating an arc emanating from the depot, it groups the nodes into several clusters, while ensuring that the total demand of each cluster is not violating the vehicle capacity. Each cluster corresponds to a TSP that can be solved by using an exact or approximate algorithm. In our experiments, we use dynamic programming to find the optimal TSP tour. After solving TSPs, the VRP solution can be obtained by combining the TSP tours. Algorithm 3 shows the pseudo-code of this algorithm.

---

**Algorithm 3** Randomized Sweep Algorithm

---

1: for each $i \in \mathcal{N}$, compute angle $\alpha_i$, respective to depot location
2: $l \leftarrow$ vehicle capacity
3: **for** $r = 1, \cdots, R$ **do**
4:     select a random angle $\alpha$
5:     $k \leftarrow 0$; initialize cluster $S_k \leftarrow \emptyset$
6:     **repeat**
7:         increase $\alpha$ until it equal to some $\alpha_i$
8:         **if** demand $d_i > l$ **then**
9:             $k \leftarrow k + 1$
10:            $S_k \leftarrow \emptyset$
11:            $l \leftarrow$ vehicle capacity
12:        **end if**
13:        $S_k \leftarrow S_k \cup \{i\}$
14:        $l \leftarrow l - d_i$
15:    **until** no unclustered node is left
16:    solve a TSP for each $S_k$
17:    merge TSP tours to produce a VRP route
18: **end for**
19: **Return**: route with the shortest length

---

# D. Asynchronous Advantage Actor-Critic for Stochastic VRP

The *Asynchronous Advantage Actor-Critic* (A3C) method proposed in (Mnih et al., 2016) is a policy gradient approach that has been shown to achieve super-human performance playing Atari games. In this paper, we utilize this algorithm for training the policy in the stochastic VRP. In this architecture, we have a central network with weights $\theta^0, \phi^0$ associated with the actor and critic, respectively. In addition, $N$ agents are running in parallel threads, each having their own set of local network parameters; we denote by $\theta^n, \phi^n$ the actor and critic weights of thread $n$. (We will use superscript $n$ to denote the operations running on thread $n$.) Each agent interacts with its own copy of the VRP at the same time as the other agents are interacting with theirs; at each time-step, the vehicle chooses the next point to visit

---

**Algorithm 4** Asynchronous Advantage Actor-Critic (A3C)

---

1: initialize the actor network with random weights $\theta^0$ and critic network with random weights $\phi^0$ in the master thread.
2: initialize $N$ thread-specific actor and critic networks with weights $\theta^n$ and $\phi^n$ associated with thread $n$.
3: **repeat**
4:     **for** each thread $n$ **do**
5:         sample a instance problem from $\Phi_{\mathcal{M}}$ with initial state $X_0^n$
6:         initialize step counter $t^n \leftarrow 0$
7:         **while** episode not finished **do**
8:             choose $y_{t+1}^n$ according to $P(y_{t+1}^n | Y_t^n, X_t^n; \theta^n)$
9:             observe new state $X_{t+1}^n$;
10:            observe one-step reward $R_t^n = R(Y_t^n, X_t^n)$
11:            let $A_t^n = \left( R_t^n + V(X_{t+1}^n; \phi) - V(X_t^n; \phi) \right)$
12:            $d\theta^0 \leftarrow d\theta^0 + \nabla_\theta A_t^n \log P(y_{t+1}^n | Y_t^n, X_t^n; \theta^n)$
13:            $d\phi^0 \leftarrow d\phi^0 + \nabla_\phi \left( A_t^n \right)^2$
14:            $t^n \leftarrow t^n + 1$
15:        **end while**
16:    **end for**
17:    periodically update $\theta^0$ using $d\theta^0$ and $\phi^0$ using $d\phi^0$
18:    $\theta^n \leftarrow \theta^0, \phi^n \leftarrow \phi^0$
19:    reset gradients: $d\theta^0 \leftarrow 0, d\phi^0 \leftarrow 0$
20: **until** training is finished

---

and receives some reward (or cost) and then goes to the next time-step. In the stochastic VRP that we consider in this paper, $R_t$ is the number of demands satisfied at time $t$. We note that the system is basically a continuous-time MDP, but in this algorithm, we consider it as a discrete-time MDP running on the times of system state changes $\{\tau_t : t = 0, \cdots\}$; for this reason, we normalize the reward $R_t$ with the duration from the previous time step, e.g., the reward is $R_t/(\tau_t - \tau_{t-1})$. The goal of each agent is to gather independent experiences from the other agents and send the gradient updates to the central network. In this approach, we periodically update the central network weights by accumulated gradients and send the updated weight to all threads. This asynchronous update procedure leads to a smooth training since the gradients are calculated from independent VRP instances.

Both actor and critic networks in this experiment are exactly the same as the ones that we employed for the classical VRP. For training the central network, we use RMSProp optimizer with learning rate $10^{-5}$.