

# Pareto Frontier Approximation Network (PA-Net) to Solve Bi-objective TSP

Ishaan Mehta, Sharareh Taghipour, and Sajad Saeedi

**Abstract**—The travelling salesperson problem (TSP) is a classic resource allocation problem used to find an optimal order of doing a set of tasks while minimizing (or maximizing) an associated objective function. It is widely used in robotics for applications such as planning and scheduling. In this work, we solve TSP for two objectives using reinforcement learning (RL). Often in multi-objective optimization problems, the associated objective functions can be conflicting in nature. In such cases, the optimality is defined in terms of Pareto optimality. A set of these Pareto optimal solutions in the objective space form a Pareto front (or frontier). Each solution has its trade-off. We present the Pareto frontier approximation network (PA-Net), a network that generates good approximations of the Pareto front for the bi-objective travelling salesperson problem (BTSP). Firstly, BTSP is converted into a constrained optimization problem. We then train our network to solve this constrained problem using the Lagrangian relaxation and policy gradient. With PA-Net we improve the performance over an existing deep RL-based method. The average improvement in the hypervolume metric, which is used to measure the optimality of the Pareto front, is 2.3%. At the same time, PA-Net has  $4.5\times$  faster inference time. Finally, we present the application of PA-Net to find optimal visiting order in a robotic navigation task/coverage planning. Our code is available on the project website<sup>1</sup>.

## I. INTRODUCTION

The travelling salesperson problem (TSP) is a popular sequencing problem. TSP generates a sequence (a.k.a. tour) that visits each city (or node) in a given graph and finally returns to the starting city. The goal of TSP is to minimize the overall traversal cost of the tour. TSP and its variants are widely used in robotics for applications like path planning for UAVs [1], multi-robot path planning [2], task allocation for robotic manipulators [3], and coverage planning [4].

We use the bi-objective travelling salesperson problem (BTSP) for application of coverage planning. Algorithms for coverage planning generate trajectories for robots to cover a given area [4]. Area coverage is used in robotic applications like cleaning robots and surveillance. Grid-based TSP planners [4] segment a given map into multiple cells and generate a coverage pattern for each cell. The optimal visiting order for these cells is generated by solving TSP that minimizes the length of the tour. We are interested in a scenario where a robot has to visit these cells and the order is dependent on two objectives. The first objective is tour length, and the second objective could be used to represent traversable condition or priority of the path. So, BTSP would be an appropriate choice in this scenario.

There are a wide variety of algorithms, ranging from exact methods to evolution-based methods that solve multi-

objective TSP [5]. Evolutionary algorithms, like non-dominated sorting genetic algorithm-II (NSGA-II) [6] and multi-objective evolutionary algorithm (MOEA/D) [7], are frequently used to tackle multi-objective TSP and other multi-objective optimization problems. Some works use evolutionary algorithms coupled with local search heuristics [8]–[10]. In practice, these evolutionary-based methods suffer in performance and computation time with an increase in the scale of the problem [11]. Another method is to use linear scalarization of the objectives using convex weights [12]. These weights are used to convert BTSP to single objective TSP, which can be solved using solvers like **OR-Tools**. The downside of the linear scalarization technique is that it is unable to find solutions in the concave region of the Pareto front [12].

**Contribution:** In this work, we address the aforementioned challenges. We present the Pareto frontier approximation network (PA-Net), a reinforcement learning (RL) based framework that generates an approximation of a set of Pareto optimal tours for BTSP. Further, we demonstrate the use of PA-Net for a robotic planning application. Our main contributions are: (1) We achieve an average improvement of 2.3% in optimality metrics along with  $4.5\times$  faster inference times as compared to the network (called DRL-MOA) proposed by Li *et al.* [13]; (2) We address the drawback of the existing approach [13], which trains separate networks to generate different solutions for approximating the Pareto front. The generalization ability of PA-Net enables it to produce a dense approximation of the Pareto front through a single network, while maintaining similar training times as DRL-MOA [13]; and (3) Our approach can be extended to generate a set of Pareto optimal solutions for other multi-objective reinforcement learning and multi-objective optimization (MOO) tasks.

**Related Work:** Sequencing problems are a subset of combinatorial optimization (CO) where the decision variables are discrete. Most CO problems are NP-Hard, and as a result state-of-the-art algorithms rely on handcrafted heuristics to make decisions that are otherwise too expensive to compute or are tailored to specific problem instances. Recently, researchers are addressing these issues using deep learning and machine learning (ML) [14]–[16].

Many recent works of deep learning-based CO methods focus on solving Euclidean TSP. Google Brain’s Pointer Network (Ptr-Net) [17] learns the conditional probability of an output sequence of elements that are discrete tokens corresponding to the positions in an input sequence. They used Ptr-Net to solve Euclidean TSP (and other CO problems) in an end-to-end fashion, where the solutions from classical methods are used as baselines for training.

Department of Mechanical and Industrial Engineering, Toronto Metropolitan University. Emails: {ishaan.mehta, sharareh, s.saeedi}@ryerson.ca

<sup>1</sup>Project website: <https://sites.google.com/view/pa-net-btsp>

Similarly, TSP was solved using a 2D graph convolution network followed by the beam search procedure [18]. Deep reinforcement learning (DRL) was used to solve various combinatorial optimization problems in [19]. Their network uses RNN-based encoder and a Ptr-Net. They trained their network using policy gradient. Deudon *et al.* [20] developed an actor-critic based architecture. They used transformers to encode the TSP graph. Further, they used 2-opt heuristics to refine the solutions generated by the network. Kool *et al.* [21] proposed a network based on attention layers. They trained their network using policy gradient. Their network outperformed other deep-learning based solvers. Further, methods like [20]–[22] have been able to close the gap (in terms of optimality) in comparison to TSP solvers like [Concorde](#) and [OR-Tools](#) by using additional heuristics or sampling procedures [20], [21].

Finding Pareto optimal solutions has been studied in deep learning literature for various multi-objective tasks. In supervised learning domain, many works focussed on multi-objective classification tasks [23]–[27]. Similarly, many multi-objective RL methods have been developed to solve multi-objective MDPs [28], [29]. Some works in RL train many single policy networks to approximate the Pareto front [13], [30]. While others have trained a single network to generate a set of Pareto optimal solutions [31], [32].

There are various methodologies to tackle MOO problems. The  $\epsilon$ -constrained methods optimize one of the objectives at a time, while using the other objectives as constraints [33], [34]. One of the most common ways is to use preference vectors or weights. These preference vectors or weights indicate the desired trade-off between various objectives. Some methods use convex weights to scalarize the objective function, and the Pareto front can be obtained by solving the optimization for multiple preference vectors [12], [35]. On the other hand, some methods use these preferences as constraints. For instance, Das *et al.* [36] generate preferences (used in constraints) based on the convex hull of individual minima of the MOO. Our method also uses preference in constraints, although in our case preferences are unit vectors sampled from the unit circle in the objective space.

Li *et al.* [13] solved multi-objective TSP by training multiple single policy networks. They converted the MOO problem into a single objective using the linear scalarization method with convex preferences. They train multiple networks with the architecture adopted from [37]. Each network is trained with a different preference weight, to approximate the Pareto front. Their network, called DRL-MOA, generated competitive results in comparison to classical methods. However, the downside of their method is that it is redundant to train multiple networks, which is time-consuming and resource-intensive. Furthermore, solutions on concave regions of the Pareto front cannot be uncovered by their network because it uses the linear scalarization technique [12]. In our work, we train a single network that can predict solutions for any unit preference vector. This enables us to produce a much denser Pareto front. Instead of using linear scalarization, our network learns to solve a

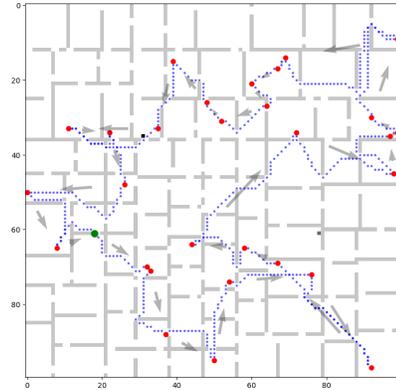


Fig. 1: A small scale BTSP tour generated by our algorithm for a map of a real-world environment. The locations of interest highlighted in red and the starting point in green. The path traversed by the robot is highlighted in blue, and arrows indicate the sequence of visiting different locations. It is clear that the path between two points of interest is not always a straight line. For such cases, we use a modified architecture of PA-Net, which can process adjacency matrices as inputs.

constrained optimization problem where the constraints are dependent on the preference vectors.

In this work, we present PA-Net, a network framework trained using the policy gradient that can approximate the Pareto front for MOO problems. Our choice of using RL is motivated by the success of deep learning-based CO methods and the fact that it is hard to generate training data for complex problems like BTSP. We use PA-Net to find a set of Pareto optimal tours for BTSP. The networks that use the policy gradient methods can be easily adopted and modified under our framework. So, the networks presented in [20], [21] are augmented for PA-Net to solve BTSP. The novelty of our algorithm is that we pose the problem of finding a set of Pareto optimal solutions as a constrained optimization problem, rather than using linear scalarization of the objective function. We use preference vectors as constraints which indicate the desired location of the solution in the objective space. Finally, we train our network using the reward constrained policy optimization, which is a paradigm used to train constrained RL applications [38]. Our network performs better than other deep learning methods in terms of quality of the Pareto front, training, and inference time. Further, we extend our framework to process non-Euclidean data. We demonstrate the use of this network for a coverage planning application.

## II. BACKGROUND

Here we review the definition of Pareto optimality and present a brief primer on solving TSP using DRL.

### A. Problem Setup

A MOO problem is defined as:

$$\min_x \vec{F}(x) = [f_1(x) f_2(x) \dots f_m(x)]^\top \quad (1)$$

where  $\vec{F}(x)$  is a vector of  $m$ -objective functions and  $x \in \mathbf{X}$  is the vector of the decision variable in  $\mathbb{R}^n$ . In such problems, often different objectives are conflicting in nature, *i.e.*, no single solution can simultaneously optimize all the objectives. Instead, a set of Pareto optimal solutions provide

the best solutions with different trade-offs between various objectives. Pareto optimality is defined as follows:

- **Dominance:** A solution  $x^a$  is said to dominate  $x^b$  ( $x^a \prec x^b$ ) if and only if  $f_i(x^a) \leq f_i(x^b)$ ,  $\forall i \in \{1, \dots, m\}$  and  $f_j(x^a) < f_j(x^b)$  such that  $\exists j \in \{1, \dots, m\}$ .
- **Pareto optimality:** A solution  $x^*$  is said to be Pareto optimal if there does not exist any solution  $x'$  such that  $x' \prec x^*$ . A set of all such points form a Pareto frontier, denoted by  $\Upsilon$ .

**Euclidean BTSP:** The Euclidean TSP is defined over a graph of  $n$  cities, where each city has coordinates  $a \in \mathbb{R}^2$ . A TSP tour  $\pi$  provides a sequence of visiting cities exactly once and then returning to the starting city. BTSP is a MOO problem that finds a set of Pareto optimal TSP tours  $\Pi$  ( $\Pi \subset \mathbb{Z}^n$ ) on a complete graph  $s$  is a sequence of  $n$  cities in a four-dimensional space  $s = [a_i^1, a_i^2]_{i=1:n}$ , where  $a_i^m \in \mathbb{R}^2$  for each  $m \in \{1, 2\}$  [13]. The goal is to find a tour  $\pi \in \Pi$  that visits each city in the graph  $s$  and can simultaneously optimize the objectives for  $m \in \{1, 2\}$ :

$$f_m(\pi) = \|a_{\pi(n)}^m - a_{\pi(1)}^m\|_2 + \sum_{i=1}^{n-1} \|a_{\pi(i)}^m - a_{\pi(i+1)}^m\|_2. \quad (2)$$

**Non-Euclidean BTSP:** In certain situations, the objective functions may not take a Euclidean form. For instance, the distance between two locations on a map might not to be a straight line because of obstacles or other path constraints, as shown in Fig. 1. So for this case, we assume that the costs between all these points are given by the adjacency matrix, *i.e.*,  $\mathbf{H}_m$  where  $m \in \{1, 2\}$ . The cost associated with the  $m^{\text{th}}$  objective for a tour  $\pi$  can be calculated using:

$$f_m(\pi) = \mathbf{H}_m[\pi(n), \pi(1)] + \sum_{i=1}^{n-1} \mathbf{H}_m[\pi(i), \pi(i+1)]. \quad (3)$$

### B. TSP using policy gradient

In our work, we adopt and modify architectures presented in [21] and [20]. These networks are trained using REINFORCE, a classic policy gradient method [39]. The input to the network is a graph  $s$ . An encoded representation of each city in the graph is obtained using an encoder. This encoded representation, along with the history of previous actions, are used by the decoder to sequentially generate a TSP tour. The actor, network  $\theta$ , is trained to minimize the total tour length given by Eq. (2) or Eq. (3). The network is trained on a batch of TSP problem instances of size  $B$ . The training objective for the actor is given by:

$$D(\theta) = \mathbb{E}_{s \sim \mathcal{S}} [\mathbb{E}_{\pi \sim p_\theta(\cdot | s)} [Q(\pi | s)]]. \quad (4)$$

Here,  $\mathcal{S}$  is the distribution from which training graphs are drawn and  $p_\theta(\pi | s)$  is the probability of a tour generated by the decoder and  $Q(\pi | s)$  is the reward which minimizes the total path length.

## III. METHODOLOGY

This section provides our mathematical formulation for MOO. We use this formulation to train a network that generates a good approximation of the Pareto front for BTSP.

### A. Problem Formulation

We intend to generate a good quality approximation of the Pareto front, which is denoted by  $\tilde{\Upsilon}$ , where  $\tilde{\Upsilon} \subset \Upsilon$ . Ideally, the set  $\tilde{\Upsilon}$  should capture a wide range of possible dominant solutions in the objective space.

BTSP is an extension of TSP to the MOO domain. Let the vector cost function for a BTSP be given by  $\vec{F}(\pi)$ , where

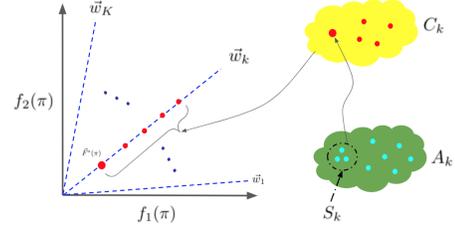


Fig. 2: Visualization of surrogate optimization of Eq. (7) for 2-D cost ( $\vec{F}(\pi)$ ) along the preference vector  $\vec{w}_k$ . All the members in the set  $C_k$  are mapped to the line along the preference vector  $\vec{w}_k$  in the objective space. The optimum point  $\vec{F}^*(\pi)$  dominates all other members in  $C_k$ . The dominant point  $\vec{F}^*$  is unique in  $C_k$ , although alternative solutions *i.e.*  $\pi_k \in S_k$  may exist that map to  $\vec{F}^*$ .

$\vec{F} \in \mathbb{R}^2$  and the tour  $\pi \in \Pi$  ( $\Pi \subset \mathbb{Z}^n$ ). The optimization problem can be written as:

$$\min_{\pi} \vec{F}(\pi) = [f_1(\pi) \ f_2(\pi)]^\top. \quad (5)$$

Where the  $i^{\text{th}}$  cost function is  $f_i : \mathbb{Z}^n \rightarrow \mathbb{R}$  for  $i \in \{1, 2\}$ . We further assume that all the cost functions are strictly positive:

$$f_i(\pi) > 0 \quad \forall i \in \{1, 2\}. \quad (6)$$

In order to find the Pareto front, we convert the MOO in Eq. (5) to a set of constrained optimization problems. This is done by discretizing the objective space using a collection of  $K$  unit preference vectors  $W : \{\vec{w}_1 \dots \vec{w}_K\}$ , where each  $\vec{w}_k \in \mathbb{R}^2$  for  $\forall k \in \{1 \dots K\}$ . These preference vectors are a set of rays emanating from the origin that uniformly divide the objective space. Each element in  $\vec{w}_k$  lies in the interval  $[0, 1]$  and  $\|\vec{w}_k\|_2 = 1$ .

The key idea is to solve a surrogate optimization problem along each preference vector in order to generate a set of dominant solutions for Eq. (5). This surrogate optimization is expressed as a set of  $K$  constrained optimization problems, where the  $k^{\text{th}}$  problem corresponding to  $\vec{w}_k \in W$  is given by:

$$\min_{\vec{F}(\pi_k)} J(\vec{F}(\pi_k)) = \|\vec{F}(\pi_k)\|_2 \quad (7)$$

*s.t.*  $\vec{F}(\pi_k) \in C_k$ .

The constraint set is defined as  $C_k = \{g(\vec{F}(\pi), \vec{w}_k) \leq 0 \mid \vec{F}(\pi) \in \mathbb{R}^2\}$  and the corresponding tour set is defined as  $A_k = \{g(\vec{F}(\pi), \vec{w}_k) \leq 0 \mid \pi \in \mathbb{Z}^n\}$ . Here, the dot product constraint  $g(\vec{F}(\pi_k), \vec{w}_k)$  is given by:

$$g(\vec{F}(\pi_k), \vec{w}_k) = 1 - \frac{\vec{w}_k \cdot \vec{F}(\pi_k)}{\|\vec{F}(\pi_k)\|_2}. \quad (8)$$

We assume  $A_k$  is non-empty. As a consequence of this assumption,  $C_k$  is also non-empty.

The constraint set  $C_k$  is a set of vector cost  $\vec{F}(\pi_k)$  where  $\pi_k \in A_k$ . Here, each  $\vec{F}(\pi_k)$  lies on the corresponding unit preference vector  $\vec{w}_k$  in objective space. The objective function in Eq. (7) minimizes the  $\mathbb{L}_2$ -norm and hence finds the points closer to origin. Below, we state a theorem that is the motivating factor of our work.

**Theorem 1.**  $\vec{F}^* \in C_k$  is the optimum solution of Eq. (7) if and only if it dominates all other points in the set.

**Proof:** Let  $\vec{F}'' \in C_k$  minimize the Eq. (7) such that  $\vec{F}^* \prec \vec{F}''$ . This dominance relation implies that  $f_i^* \leq f_i'' \forall i \in \{1, 2\}$  and  $\exists j \in \{1, 2\}$  such that  $f_j^* < f_j''$ . This dominance relation leads to the following result:

$$\|\vec{F}^*\|_2 < \|\vec{F}''\|_2. \quad (9)$$

But this result is a contradiction because  $\vec{F}''$  minimizes Eq. (7). Hence,  $\vec{F}^* \in C_k$  that dominates all other points in the set are the optimum solution for problem Eq. (7). ■

A similar argument can be made to show that a non-dominated solution in  $C_k$  is the optimum solution of Eq. (7). The dominant point in the set, *i.e.*,  $\vec{F}^* \in C_k$  is also unique in the set. An intuitive proof for this can be visualized using the case for  $C_k \subset \mathbb{R}^2$  as shown in Fig. 2. Because of the dot product constraint, all the possible members in the set lie on the unit vector  $\vec{w}_k$ . It is clear from Fig. 2 that the point in the set  $C_k$  closest to origin dominates all other points and is, in fact, the optimum solution of Eq. (7). Further, there can be multiple solutions in  $A_k$  that lead to the dominant objective value, *i.e.*,  $\vec{F}^*$ . Mathematically, the solution set  $S_k \subset A_k$ , where  $S_k = \{\vec{F}(\pi_k) = \vec{F}^* \mid \pi_k \in A_k\}$ , all members of  $S_k$  will generate dominant objective values.

Essentially, Theorem 1 demonstrates the viability of approximating the Pareto front for the problem in Eq. (5) through the surrogate optimization problem in Eq. (7). Solving the optimization problem in Eq. (7) for large values of  $K$  can be computationally intractable. We address this issue through the generalization power of deep neural networks. PA-Net learns to approximately solve Eq. (7) on a given input preference set. We select a sparse preference set (more details in Sec. III-B) that captures the diversity of the solutions in the objective space. This enables PA-Net to generalize to unseen preferences as well. Unlike linear scalarization methods, our formulation can also find concave Pareto frontier. We demonstrate this with an example of a concave Pareto front in the Appendix.

### B. Network Architecture

**Euclidean BTSP:** Our training framework can easily be used by any existing policy gradient-based networks to solve TSP. In this work, we adopt and modify the architecture from existing works to solve Euclidean BTSP. The first network we use is Encode Attend Navigate (EAN) network [20]. The second network we use is the Attention network (AT) from [21]. We refer to these modified networks trained with our framework as **PA-EAN** and **PA-AT**, respectively.

The input for Euclidean BTSP is 4D coordinates of the cities, where each pair of coordinates is used to calculate the corresponding objective. To generate a set of dominant tours, we augment these network architecture by adding an input of a set of preference vectors  $W$  of size  $K$ . Each  $\vec{w}_k \in W$  is encoded in higher dimensions using a single layer feed-forward network. These additional layers learn features corresponding to different preferences. This encoding is combined with the encoded representation of the graph and then passed on to the decoder. With this architecture, the network can be trained for various preferences.

**Non-Euclidean BTSP:** For the case where the input is an adjacency matrix, we are required to learn a representation of each city corresponding to each objective. This representation is used as an input to the network to generate a set of TSP tours. A graph transformer-based encoder by Sykora *et al.* [40] is used to learn this representation. The input to

the encoder is the adjacency matrix  $\mathbf{H}_j$  and an initial set of features of each city in the graph,  $s$  given by  $s^j = [a_1^j \dots a_n^j]$  where  $j \in \{1, 2\}$ . The description of the initial feature for the  $i^{\text{th}}$  city corresponding to the  $j^{\text{th}}$  objective,  $a_i^j \in \mathbb{R}^3$ , is given in Table I.

TABLE I: Description of initial city features for PA-AD

Name	Dim.	Type
Sum of neighbouring edge weights	1	float
Min. neighbouring edge weight	1	float
Max. neighbouring edge weight	1	float

The graph encoder produces an encoded representation of the features for each objective, *i.e.*,  $\vec{s}^j$ . It should be noted that we process features for each objective independent of the other. Finally, we stack the learned features for each objective to obtain a combined representation  $\vec{\mathbf{S}} = [\vec{s}^1 \vec{s}^2]^\top$ . Now,  $\vec{\mathbf{S}}$  along with the preference encoding, is used by the network to predict BTSP tours. For this case, we use the AT network [21] along with the preference encoder. We will refer to this combined network as **PA-AD**.

### C. Training Methodology

The reward-constrained policy optimization is an actor-critic algorithm that solves constrained RL problems [38]. It uses a Lagrangian of the constrained problem as the objective function, where after each gradient update step, the Lagrangian multipliers are updated based on the constraint violation. We use the reward-constrained policy optimization to train a network to solve the problem in Eq. (7) for all  $k \in \{1 \dots K\}$ .

We intend to train a single network that generates a set of dominant tours  $T : \{\pi_1, \dots, \pi_K\}$ . Hence, the problem in Eq. (7) for each  $\pi_k \in T$  can be written in the parametric format as:

$$\begin{aligned} \min_{\theta} J(\pi_k(\theta, s_j)) &= \|\vec{F}(\pi_k(\theta, s_j))\|_2 \\ \text{s.t. } g_k(\vec{F}(\pi_k(\theta, s_j), \vec{w}_k)) &\leq 0. \end{aligned} \quad (10)$$

Here,  $\theta$  is the parameters of the actor network and  $\pi_k(\theta, s_j)$  is the tour generated by the actor network corresponding to  $k^{\text{th}}$  preference for the input graph  $s_j$ . For notational convenience, we denote  $\pi_k(\theta, s_j)$  as  $\pi_k^j$ . The Lagrangian dual problem for Eq. (10) is:

$$L_k(\pi_k^j, \lambda_k) = \max_{\lambda_k \geq 0} \min_{\theta} J(\pi_k^j) + \lambda_k \cdot g_k(\vec{F}(\pi_k^j, \vec{w}_k)). \quad (11)$$

Here,  $\lambda_k$  is the  $k^{\text{th}}$  Lagrangian multiplier corresponding to the preference vector  $\vec{w}_k$ . We use the Lagrangian in Eq. (11) as the reward for the network. Based on this reward, the training objective for the actor can be written in our case as:

$$D_{AC}(\theta) = \mathbb{E}_{s_j \sim S} [\mathbb{E}_{\vec{w}_k \sim W} [\mathbb{E}_{\pi \sim p_{\theta}(\cdot | s_j)} [L_k(\pi_k^j, \lambda_k)]]]. \quad (12)$$

A critic network is used to provide predictions  $b_{\phi}(\vec{w}_k, s_j)$  on the reward given in Eq. (11). This critic network is trained on the mean squared error between its predictions and rewards of the actor, which is given by:

$$D_{CR}(\phi) = \frac{1}{K \cdot B} \sum_{k=1}^K \sum_{j=1}^B \|b_{\phi}(\vec{w}_k, s_j) - L_k(\pi_k^j, \lambda_k)\|_2^2. \quad (13)$$

The gradient for the training of the actor network is approximated using REINFORCE [39]:

$$\begin{aligned} \nabla_{\theta} D_{AC}(\theta) &\approx \frac{1}{K \cdot B} \sum_{k=1}^K \sum_{j=1}^B [L_k(\pi_k^j, \lambda_k) - \\ &b_{\phi}(\vec{w}_k, s_j)] \cdot \nabla_{\theta} \log(p_{\theta}(\pi_k^j)). \end{aligned} \quad (14)$$

The description of the training of PA-Net is given in **Algorithm 1**. We start with the initialization of weights and learning rates for the network and the set of preference vectors, along with other hyperparameters that are the ascent rate of the Lagrangian multipliers  $\alpha$  and  $[\lambda_{min}, \lambda_{max}]$  the limits for the multipliers. The network is trained for  $N$  iterations. At each iteration, a batch of graphs  $\Omega$  of size  $B$  is generated. For each  $s_j \in \Omega$  and the corresponding preference vector,  $\vec{w}_k \in W$  a tour  $\pi_k^j$  is constructed by the network. Based on generated tours, the objective for actor and critic are calculated. Then parameters of the network are updated using gradient descent. At the end of each iteration, the Lagrangian multiplier corresponding to each preference vector is updated in an ascent step using:

$$\lambda_k^{i+1} = \Gamma_\lambda(\lambda_k^i + \frac{\alpha}{B} \cdot \sum_{j=1}^B g_k(F(\pi_k^j, \vec{w}_k))). \quad (15)$$

Here  $\Gamma_\lambda(\cdot)$  ensures that the multipliers remain within the limits, *i.e.*,  $[\lambda_{min}, \lambda_{max}]$  and  $\alpha$  is prespecified ascent rate. Although the network is trained on a fixed set of preferences  $W$ , it can generalize to a larger set of preferences. For training of BTSP, we generate preferences by sampling the unit circle in  $\mathbb{R}^2$ . Since the objective functions are strictly positive, the angle for unit vectors is sampled from the interval  $(0^\circ, 90^\circ)$ . A higher number of preferences, *i.e.*, would lead to better performance. But increasing  $K$  also increases training times. So, to keep the training tractable, we limit the number to  $K = 20$ .

One downside of our approach is that our network might be susceptible to local minima. However, in practice, our approach can generate competitive results as presented next.

---

### Algorithm 1: Training of PA-Net

---

**input :**  $[\theta, \phi, \eta_A, \eta_C], [W, \alpha, \lambda_{1:K}, \lambda_{min}, \lambda_{max}] \leftarrow$  Initialization of network weights, set of preference vectors and corresponding parameters.

**output:** Trained network parameters of PA-Net  $\theta^*, \phi^*$ .

```

for  $i \leftarrow 1, 2, \dots, N$  do
   $\Omega: \{s_1 \dots s_B\} \leftarrow$  Sample a Batch of TSP Graphs of size  $B$  from distribution  $\mathcal{S}$ .
  for  $k \leftarrow 1, 2, \dots, K$  do
    for  $j \leftarrow 1, 2, \dots, B$  do
       $\pi_k^j \leftarrow$  Actor network Generates TSP Tour for each  $s_j$  and  $\vec{w}_k$ .
       $b_\phi(\vec{w}_k, s_j) \leftarrow$  Critic Network predicts the baseline
       $L_k(\pi_k^j, \lambda_k) \leftarrow$  Calculate the Lagrangian using Eq. (11).
    Actor Update:  $\theta \leftarrow \theta - \eta_A \cdot \nabla_\theta D_{AC}(\theta)$ 
    Critic Update:  $\phi \leftarrow \phi - \eta_C \cdot \nabla_\phi D_{cr}(\phi)$ 
  for  $k \leftarrow 1, 2, \dots, K$  do
     $\lambda_k \leftarrow$  Update the Lagrangian multipliers using Eq. (15)

```

---

## IV. EXPERIMENTS

To evaluate the efficacy of PA-Net, we present two experiments. The first experiment is on BTSP instances, where the input is Euclidean data. The second experiment is an application for coverage planning, where we test our network on non-Euclidean input data generated from simulated grid world maps.

The performance of PA-Net is compared with deep learning-based method DRL-MOA [13], evolution-based strategies: NSGA-II and MOEA/D [7]. Further, we generate baseline results using TSP solver from **OR-Tools** library.

Note that in this case, we use linear scalarization of objectives. The experiments for deep learning methods are carried out on NVIDIA V100 Volta GPU. Whereas, for NSGA-II, MOEA/D and OR-Tools experiments are carried out on dual-core Intel i5 processor.

The results of PA-Net are compared with other methods on the basis of hypervolume (HV) of the solutions. HV is a hybrid metric used to evaluate Pareto fronts [41], [42]. It represents the volume covered by the non-dominated set of solutions with respect to the worst-case solution. A higher HV indicates a better quality of the Pareto front, both in terms of optimality and coverage of the objective space. Here, HV is calculated as the percentage of points dominated by the solution set out of densely sampled points in a fixed volume in the objective space. This volume is generated with respect to a reference point. We use a fixed reference point to compute HV for all the algorithms. It is calculated as the product of the number of cities and unity vector, for example, reference point for 200-city BTSP is  $[200.0, 200.0]^T$ .

We also mention the run time of each algorithm. However, evolutionary-based methods and OR-Tools are CPU-based methods, whereas DRL-MOA and PA-Net are GPU-based methods. Libraries like OR-tools are highly optimized and can generate solutions in a short time. On the other hand, deep-learning based methods have relatively higher run-time, but they can take advantage of parallelization and generate a batch of solutions in one shot.

TABLE II: Training details for our networks and DRL-MOA. Our network can learn to generalize for various preferences for similar training time to DRL-MOA

	PA-EAN	PA-AT	PA-AD	DRL-MOA
Batch Size	60	60	60	200
Epochs	1	3	4	1/net.
Steps (per epoch)	20000	5000	5000	2500
Input Graph Size	$120 \times 4$	$40 \times 4$	$40 \times 40$	$40 \times 4$
Training Times (hrs)	$\sim 12$	$\sim 8$	$\sim 10$	$\sim 9$

**Training Setup:** For Euclidean BTSP, we train two architectures, *i.e.*, PA-AT and PA-EAN. For the non-Euclidean BTSP, we train PA-AD. Training details of these networks along with DRL-MOA are given in Table IV. The training times are reported based on training on NVIDIA V100 Volta GPU. All the networks except PA-EAN are trained on 40-city BTSP instances. Further, all our networks are trained for  $K = 20$  preferences. PA-AD network is trained on adjacency matrices.

### A. Evaluation on Euclidean BTSP Instances

In this experiment, we evaluate the performance of PA-Net on bi-objective TSP. The input graph  $G$  is a sequence of  $n$  cities, where each city  $a_i \in \mathbb{R}^4$  is represented by 4 dimensional features. We use  $\mathbb{L}_2$  norm using Eq. (2) for the two objectives. For each algorithm, we generate 100 solutions corresponding to various preferences among the two objectives. It should be noted that the preferences used in linear scalarization are sampled from  $[0,1]$  and satisfy convexity constraints. In PA-Net, the preferences used are unit vectors that uniformly segment the objective space.

A set of 25 synthetically generated BTSP instances are used in this experiment. The average HV and run times for each algorithm are reported in Table III. NSGA-II and

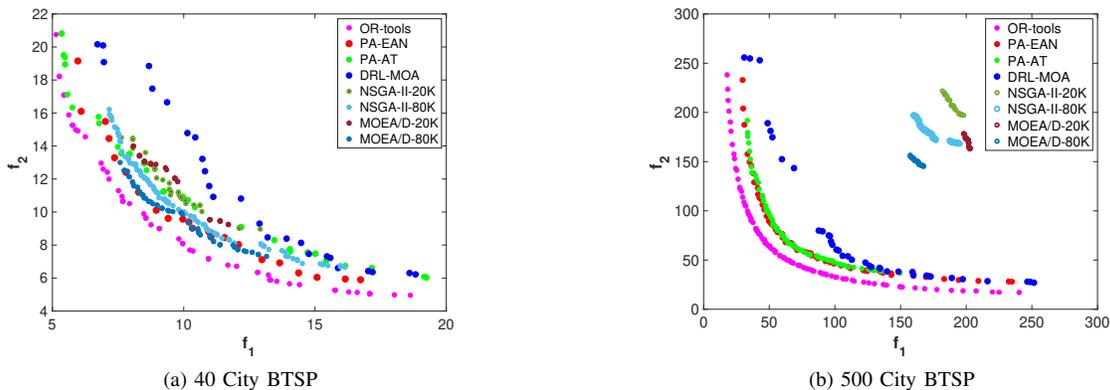


Fig. 3: Visualization of the dominant solutions for different problem instances. It can be seen that our network (PA-EAN and PA-AT) generates significantly better objective values than DRL-MOA and evolutionary methods

TABLE III: Quantitative comparison of solutions for BTSP. Our methods outperform DRL-MOA and other evolutionary algorithms in terms of HV and time.

Algo.	40-City		200-City		500-City		1000-City	
	HV (%)	Time (s)	HV (%)	Time (s)	HV (%)	Time (s)	HV (%)	Time (s)
NSGA-II (20K)	67.3	5.64	45.4	8.04	38.4	14.7	33.8	27.1
NSGA-II (80K)	72.5	21.7	53.9	30.8	46.36	58.7	41.48	107.3
MOEA/D (20K)	66.7	9.357	47.5	12.7	40.4	20.7	35.7	33.5
MOEA/D (80K)	70.7	34.65	55.98	48.2	48.8	79.53	43.89	130.75
DRL-MOA	73.6	5.87	80.63	29.3	84.5	72.9	85.9	145.5
PA-EAN (ours)	75.4	<b>1.59</b>	83.2	<b>6.03</b>	86.92	<b>15.14</b>	88.45	<b>30.38</b>
PA-AT (ours)	74.18	2.6	80.9	13.4	85.15	33.5	87.35	68.2
OR-tools ( $n_l = \text{NO LIMIT}$ )	<b>78.14</b>	2.16	<b>86.5</b>	86.8	<b>91.07</b>	732	<b>93.39</b>	3730
OR-tools ( $n_l = 10$ )	78.08	0.93	84.8	10.4	89.8	53.5	92.38	209

TABLE IV: Quantitative comparison of Pareto front for Coverage Planning. Our algorithm generates promising initial results.

Algo.	40-City		100-City		200-City	
	HV (%)	Time (s)	HV (%)	Time (s)	HV (%)	Time (s)
PA-AD (ours)	71.59	3.06	72.14	6.26	76.4	11.85
OR-tools ( $n_l = 100$ )	75.9	2.15	78.7	17.4	84.29	82.5
OR-tools ( $n_l = 10$ )	66.6	0.75	68.5	2.72	82.4	9.9

MOEA/D are tested at different values for the maximum number of iterations. For OR-Tools, we report results for different values of the maximum number of iterations for solution refinement,  $n_l$ , where the solver terminates once the specified limit is reached. The results of the quantitative comparison for BTSP is given in Table III. Visualization of dominant objective values attained by different algorithms is shown in Fig. 3 (a)-(b).

**Discussion:** It is clear that the baseline method *i.e.*, OR-Tools achieves the best values for HV metric. Further, these values improve with increasing the maximum solution limit of the solver. The routing problem solver of OR-Tools is a highly optimized library built on years of research by the operation research community. It should be mentioned, that even the state-of-the-art deep learning-based methods to solve single objective TSP are able to compete against OR-Tools by using additional neural local heuristics or sampling a large set of solution space [20]–[22]. On the other hand, evolutionary methods significantly underperform as the scale of the problem increases. This is likely because these algorithms are unable to explore the solution space well for larger problems. Running these algorithms for more iterations could potentially improve their performance in terms of HV, but this comes with an additional computational cost. All deep learning-based methods achieve competitive results in terms of HV. Our networks achieve much better performance as compared to DRL-MOA in terms of HV. The

average improvement in HV over DRL-MOA for PA-EAN and PA-AT is 2.3% and, 0.7% respectively. Further, PA-EAN and PA-AT generates the complete Pareto front about  $4.5\times$  and  $2.1\times$  faster, respectively, as compared to DRL-MOA. Further, we have either lower or comparable training times relative to DRL-MOA, see Table II. For each problem set, PA-Net can infer a solution from a single network, whereas DRL-MOA has to train and rely on multiple networks. Another notable point is that DRL-MOA has many gaps in its Pareto front, see Fig. 3. This demonstrates the ineffectiveness of linear scalarization approach in deep learning setting.

### B. Application for Coverage Planning

We use BTSP for the application of coverage planning. In this experiment, the robot has to visit various locations in a given map while optimizing for total distance travelled and an additional priority metric. An example of such a map is shown in Fig. 1. Such a priority metric can be representative of attributes like traversable conditions of the path or traffic on a given path *etc.* The distances between locations are stored in adjacency matrices for this case because the distances are not Euclidean.

For this experiment, we synthetically generate 10 maps using PathBench [43], an open-source motion-planning benchmarking platform. We then randomly sample points from the free space in the map as locations of interest, as shown in Fig. 1. We assume that the graph formed by these

points are fully connected. The input adjacency matrix for the first objective,  $\mathbf{H}_1$ , is generated by determining path-lengths between all the points using A\* algorithm. For the second objective, we randomly sample  $2D$  points where each coordinate is in range  $[0,1)$ . We then generate the second adjacency matrix  $\mathbf{H}_2$ , by computing the Euclidean distance between these points. For this dataset, we test the performance of our PA-AD network while using OR-Tools as the baseline. The average results for HV and runtimes are reported in Table IV.

**Discussion:** OR-Tools finds much better solutions when this limit is increased. As mentioned before, OR-Tools is a highly optimized solver for TSP and other routing problems. Like the Euclidean case, there is a lag in the performance of our network as compared to OR-Tools. This is likely because the network converges to local minima. Nonetheless, these initial results look encouraging as we are able to outperform existing deep learning approach, *i.e.*, DRLMOA.

## V. CONCLUSIONS

We presented PA-Net, a network that approximates the Pareto frontier for the bi-objective TSP. Our results indicate a competitive performance in terms of optimality of the solutions. This is achieved by segmenting the objective space using a set of unit vectors which represent trade-offs among various objectives. We then use these preference vectors to convert the unconstrained optimization problem into a set of constrained optimization problems. Then the network is trained using Lagrangian relaxation and policy gradient to generate solutions for these constrained problems. While PA-Net is trained on a fixed number of preference vectors, it generalizes well to other unseen preferences as well. The effectiveness of our method is highlighted by the significant gains made in terms of quality of solutions, inference and training times. Although we focus on bi-objective TSP in this work, our training framework can be applied to other MOO problems. We also demonstrated a use case of PA-Net for a coverage planning application. Our future investigation would be focussed on bridging the performance gap with OR-Tools through informed use of local heuristics and tackling the issue of convergence to local minima by incorporating expert training data from solvers like OR-Tools and Concorde. Further, we also plan on extending this work to cases with dynamic costs and multi-robot systems.

## APPENDIX

### Ablation Studies

In this ablation study, we evaluate the contribution of the preference layer. We train two networks, *i.e.*, PA-EAN<sup>-</sup> and PA-AT<sup>-</sup>. The training conditions for both these networks are kept the same as our proposed networks, with the only difference being that the preference embedding layer is removed for both these networks. We compare our proposed network with these on the basis of hypervolume. The results are reported in Table V. It can be clearly seen that without the preference layer, the performance of the networks drop significantly.

In the case of non-Euclidean BTSP, a key addition is the graph transformer layers that learn representation for

TABLE V: Impact of using preference encoder. It can be seen that preference encoders play a critical role for the performance of our networks.

Network	40-City	200-City	500-City	1000-City
	HV(%)	HV(%)	HV(%)	HV(%)
PA-EAN	75.4	83.2	86.92	88.4
PA-EAN <sup>-</sup>	39.1	25.6	25	24.6
PA-AT	74.18	80.9	85.15	87.3
PA-AT <sup>-</sup>	52	48.9	52.9	55.8

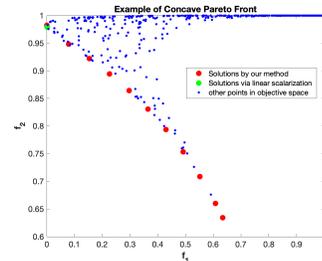


Fig. 4: The above plot depicts the convergence of our method to Concave Pareto Front. Here red points are the results generated by our algorithm, green points are the solution from linear scalarization method and blue points represent the set of possible solutions in the objective space.

initial features. In order to evaluate its impact, we train a network, *i.e.*, PA-AD<sup>-</sup> without a graph convolution layer. In this case, we use the initial feature set described in Table I along with the coordinates of the city as the input to the network. We then compare PA-AD and PA-AD<sup>-</sup> on the map dataset generated for the non-Euclidean BTSP experiment. The results for comparison are reported in Table VI. It can be clearly seen that graph transformer encoder plays a crucial role to learn useful representations of features for each city.

TABLE VI: Impact of using graph transformer encoder. Using a graph encoder enables the network to learn better features as indicated by higher HV values.

Network	40-City	100-City	200-City
PA-AD	71.59	72.14	76.4
PA-AD <sup>-</sup>	60	53	53.2

### Convergence to Concave Pareto Fronts

We solve a concave MOO problem to demonstrate the convergence of our optimization framework given by Eq. (7). The concave problem taken from [24] is given by:

$$\min_x \vec{F}(x) = [f_1(x), f_2(x)]^\top, \quad (16)$$

where,

$$f_1(x) = 1 - \exp(-\sum_{i=1}^d (x_i - \frac{1}{\sqrt{d}})^2), \quad (17)$$

$$f_2(x) = 1 - \exp(-\sum_{i=1}^d (x_i + \frac{1}{\sqrt{d}})^2).$$

Here  $x = [x_1, x_2]^\top \in \mathbb{R}^{2+}$  and  $d = 2$ . The surrogate optimization in this case with preference  $\vec{w}_k$  is given by:

$$\min_{\vec{F}(x^k)} J(\vec{F}(x^k)) = \|\vec{F}(x^k)\|_2$$

$$s.t. 1 - \frac{\vec{w}_k \cdot \vec{F}(x^k)}{J(\vec{F}(x^k))} \leq 0 \quad (18)$$

We solve the above problem in Matlab for  $K = 20$ . The preference is generated by  $\vec{w}_k = [\cos(\phi_k), \sin(\phi_k)]^\top$ , where  $\phi_k \in (0, 90)$ . We also solve the MOO problem with a simple linear scalarization of objective. In this case, the preference is given by  $\alpha^k = [\alpha_1, \alpha_2]^\top \in \mathbb{R}^{2+}$  such that  $\alpha_1 + \alpha_2 = 1$ . We use  $K = 100$  preferences in this case. The  $k^{th}$  objective function for linear scalarization is:

$$\min_{\vec{F}(x^k)} R(\vec{F}(x^k)) = \alpha_1 \cdot f_1(x^k) + \alpha_2 \cdot f_2(x^k) \quad (19)$$

The results are shown in Fig. 4. It can be clearly seen that our method is able to produce the concave Pareto front. On the other hand, linear scalarization is unable to find solutions

on the concave part of the Pareto front. This example demonstrates that our method can certainly be extended to MOO with concave Pareto fronts.

## REFERENCES

- [1] Y. Xu and C. Che, "A brief review of the intelligent algorithm for traveling salesman problem in UAV route planning," in 2019 IEEE 9th International Conference on Electronics Information and Emergency Communication (ICEIEC). IEEE, 2019, pp. 1–7.
- [2] Z. Yu, L. Jinhai, G. Guochang, Z. Rubo, and Y. Haiyan, "An implementation of evolutionary computation for path planning of cooperative mobile robots," in Proceedings of the 4th World Congress on Intelligent Control and Automation (Cat. No. 02EX527), vol. 3. IEEE, 2002, pp. 1798–1802.
- [3] P. T. Zacharia and N. Aspragathos, "Optimal robot task scheduling based on genetic algorithms," Robotics and Computer-Integrated Manufacturing, vol. 21, no. 1, pp. 67–79, 2005.
- [4] R. Bormann, F. Jordan, J. Hampp, and M. Hägele, "Indoor coverage path planning: Survey, implementation, analysis," in 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 1718–1725.
- [5] T. Lust and J. Teghem, "The multiobjective traveling salesman problem: a survey and a new approach," in Advances in Multi-Objective Nature Inspired Computing. Springer, 2010, pp. 119–141.
- [6] B. A. Beirigo and A. G. dos Santos, "Application of NSGA-II framework to the travel planning problem using real-world travel data," in 2016 IEEE Congress on Evolutionary Computation (CEC). IEEE, 2016, pp. 746–753.
- [7] W. Peng, Q. Zhang, and H. Li, "Comparison between moea/d and NSGA-II on the multi-objective travelling salesman problem," in Multi-objective memetic algorithms. Springer, 2009, pp. 309–324.
- [8] A. Jaskiewicz, "On the performance of multiple-objective genetic local search on the 0/1 knapsack problem—a comparative experiment," IEEE Transactions on Evolutionary Computation, vol. 6, no. 4, pp. 402–412, 2002.
- [9] L. Ke, Q. Zhang, and R. Battiti, "Hybridization of decomposition and local search for multiobjective optimization," IEEE transactions on cybernetics, vol. 44, no. 10, pp. 1808–1820, 2014.
- [10] X. Cai, Y. Li, Z. Fan, and Q. Zhang, "An external archive guided multi-objective evolutionary algorithm based on decomposition for combinatorial optimization," IEEE Transactions on Evolutionary Computation, vol. 19, no. 4, pp. 508–523, 2014.
- [11] X. Zhang, Y. Tian, R. Cheng, and Y. Jin, "A decision variable clustering-based evolutionary algorithm for large-scale many-objective optimization," IEEE Transactions on Evolutionary Computation, vol. 22, no. 1, pp. 97–112, 2016.
- [12] S. Boyd, S. P. Boyd, and L. Vandenberghe, Convex optimization. Cambridge university press, 2004.
- [13] K. Li, T. Zhang, and R. Wang, "Deep reinforcement learning for multiobjective optimization," IEEE Transactions on Cybernetics, 2020.
- [14] Y. Bengio, A. Lodi, and A. Prouvost, "Machine learning for combinatorial optimization: a methodological tour d’horizon," European Journal of Operational Research, 2020.
- [15] N. Vesselinova, R. Steinert, D. F. Perez-Ramirez, and M. Boman, "Learning combinatorial optimization on graphs: A survey with applications to networking," IEEE Access, vol. 8, pp. 120 388–120 416, 2020.
- [16] N. Mazyavkina, S. Sviridov, S. Ivanov, and E. Burnaev, "Reinforcement learning for combinatorial optimization: A survey," arXiv preprint arXiv:2003.03600, 2020.
- [17] O. Vinyals, M. Fortunato, and N. Jaitly, "Pointer networks," in Advances in neural information processing systems, 2015, pp. 2692–2700.
- [18] C. K. Joshi, T. Laurent, and X. Bresson, "An efficient graph convolutional network technique for the travelling salesman problem," arXiv preprint arXiv:1906.01227, 2019.
- [19] I. Bello, H. Pham, Q. V. Le, M. Norouzi, and S. Bengio, "Neural combinatorial optimization with reinforcement learning," arXiv preprint arXiv:1611.09940, 2016.
- [20] M. Deudon, P. Cournut, A. Lacoste, Y. Adulyasak, and L.-M. Rousseau, "Learning heuristics for the tsp by policy gradient," in International conference on the integration of constraint programming, artificial intelligence, and operations research. Springer, 2018, pp. 170–181.
- [21] W. Kool, H. Van Hoof, and M. Welling, "Attention, learn to solve routing problems!" arXiv preprint arXiv:1803.08475, 2018.
- [22] Y. Ma, J. Li, Z. Cao, W. Song, L. Zhang, Z. Chen, and J. Tang, "Learning to iteratively solve routing problems with dual-aspect collaborative transformer," Advances in Neural Information Processing Systems, vol. 34, 2021.
- [23] O. Sener and V. Koltun, "Multi-task learning as multi-objective optimization," in Advances in Neural Information Processing Systems, 2018, pp. 527–538.
- [24] X. Lin, H.-L. Zhen, Z. Li, Q.-F. Zhang, and S. Kwong, "Pareto multi-task learning," in Advances in Neural Information Processing Systems, 2019, pp. 12 060–12 070.
- [25] D. Mahapatra and V. Rajan, "Multi-task learning with user preferences: Gradient descent with controlled ascent in pareto optimization," in International Conference on Machine Learning. PMLR, 2020, pp. 6597–6607.
- [26] M. Ruchte and J. Grabocka, "Efficient multi-objective optimization for deep learning," arXiv preprint arXiv:2103.13392, 2021.
- [27] A. Navon, A. Shamsian, G. Chechik, and E. Fetaya, "Learning the pareto front with hypernetworks," in International Conference on Learning Representations, 2021. [Online]. Available: <https://openreview.net/forum?id=NjF772F4ZZR>
- [28] D. M. Roijers, P. Vamplew, S. Whiteson, and R. Dazeley, "A survey of multi-objective sequential decision-making," Journal of Artificial Intelligence Research, vol. 48, pp. 67–113, 2013.
- [29] S. Parisi, M. Pirodda, N. Smacchia, L. Bascetta, and M. Restelli, "Policy gradient approaches for multi-objective sequential decision making," in 2014 International Joint Conference on Neural Networks (IJCNN). IEEE, 2014, pp. 2323–2330.
- [30] P. Vamplew, R. Issabekov, R. Dazeley, C. Foale, A. Berry, T. Moore, and D. Creighton, "Steering approaches to pareto-optimal multiobjective reinforcement learning," Neurocomputing, vol. 263, pp. 26–38, 2017.
- [31] R. Yang, X. Sun, and K. Narasimhan, "A generalized algorithm for multi-objective reinforcement learning and policy adaptation," in Advances in Neural Information Processing Systems, 2019, pp. 14 636–14 647.
- [32] S. Parisi, M. Pirodda, and M. Restelli, "Multi-objective reinforcement learning through continuous pareto manifold approximation," Journal of Artificial Intelligence Research, vol. 57, pp. 187–227, 2016.
- [33] G. Mavrotas, "Effective implementation of the  $\epsilon$ -constraint method in multi-objective mathematical programming problems," Applied mathematics and computation, vol. 213, no. 2, pp. 455–465, 2009.
- [34] A. Chinchuluun and P. M. Pardalos, "A survey of recent developments in multiobjective optimization," Annals of Operations Research, vol. 154, no. 1, pp. 29–50, 2007.
- [35] C. C. Coello, C. Dhaenens, and L. Jourdan, Advances in multi-objective nature inspired computing. Springer, 2009, vol. 272.
- [36] I. Das and J. E. Dennis, "Normal-boundary intersection: A new method for generating the pareto surface in nonlinear multicriteria optimization problems," SIAM journal on optimization, vol. 8, no. 3, pp. 631–657, 1998.
- [37] M. Nazari, A. Oroojlooy, L. Snyder, and M. Takác, "Reinforcement learning for solving the vehicle routing problem," Advances in neural information processing systems, vol. 31, 2018.
- [38] C. Tessler, D. J. Mankowitz, and S. Mannor, "Reward constrained policy optimization," arXiv preprint arXiv:1805.11074, 2018.
- [39] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," Machine learning, vol. 8, no. 3–4, pp. 229–256, 1992.
- [40] Q. Sykora, M. Ren, and R. Urtasun, "Multi-agent routing value iteration network," in International Conference on Machine Learning. PMLR, 2020, pp. 9300–9310.
- [41] J.-F. Bérubé, M. Gendreau, and J.-Y. Potvin, "An exact phi-constraint method for bi-objective combinatorial optimization problems: Application to the traveling salesman problem with profits," European journal of operational research, vol. 194, no. 1, pp. 39–50, 2009.
- [42] C. Audet, J. Bigeon, D. Cartier, S. Le Digabel, and L. Salomon, "Performance indicators in multiobjective optimization," European journal of operational research, 2020.
- [43] A.-I. Toma, H.-Y. Hsueh, H. A. Jaafar, R. Murai, P. H. Kelly, and S. Saeedi, "PathBench: A benchmarking platform for classical and learned path planning algorithms," in 2021 18th Conference on Robots and Vision (CRV). IEEE, 2021, pp. 79–86.