

SCHEDULENET: LEARN TO SOLVE MINMAX MULTIPLE TRAVELLING SALESMAN PROBLEM

Anonymous authors

Paper under double-blind review

ABSTRACT

There has been continuous effort to learn to solve famous CO problems such as Traveling Salesman Problem (TSP) and Vehicle Routing Problem (VRP) using reinforcement learning (RL). Although they have shown good optimality and computational efficiency, these approaches have been limited to scheduling a single-vehicle. *MinMax* mTSP, the focus of this study, is the problem seeking to minimize the total completion time for multiple workers to complete the geographically distributed tasks. Solving *MinMax* mTSP using RL raises significant challenges because one needs to train a distributed scheduling policy inducing the cooperative strategic routings using only the single delayed and sparse reward signal (makespan). In this study, we propose the ScheduleNet that can solve mTSP with any numbers of salesmen and cities. The ScheduleNet presents a state (partial solution to mTSP) as a set of graphs and employs type aware graph node embeddings for deriving the cooperative and transferable scheduling policy. Additionally, to effectively train the ScheduleNet with sparse and delayed reward (makespan), we propose an RL training scheme, Clipped REINFORCE with "target net," which significantly stabilizes the training and improves the generalization performance. We have empirically shown that the proposed method achieves the performance comparable to Google OR-Tools, a highly optimized meta-heuristic baseline.

1 INTRODUCTION

There have been numerous approaches to solve combinatorial optimization (CO) problems using machine learning. Bengio et al. (2020) have categorized these approaches into *demonstration* and *experience*. In *demonstration* setting, supervised learning has been employed to mimic the behavior of the existing expert (e.g., exact solvers or heuristics). On the other hand, in the *experience* setting, typically, reinforcement learning (RL) has been employed to learn a parameterized policy that can solve newly given target problems without direct supervision. While the demonstration policy cannot outperform its guiding expert, RL-based policy can outperform the expert because it improves its policy using a reward signal. Concurrently, Mazyavkina et al. (2020) have further categorized the RL approaches into *improvement* and *construction* heuristics. An *improvement* heuristics start from the arbitrary (complete) solution of the CO problem and iteratively improve it with the learned policy until the improvement stops (Chen & Tian, 2019; Ahn et al., 2019). On the other hand, the *construction* heuristics start from the empty solution and incrementally extend the partial solution using a learned sequential decision-making policy until it becomes complete.

There has been continuous effort to learn to solve famous CO problems such as Traveling Salesman Problem (TSP) and Vehicle Routing Problem (VRP) using RL-based construction heuristics (Bello et al., 2016; Kool et al., 2018; Khalil et al., 2017; Nazari et al., 2018). Although they have shown good optimality and computational efficiency performance, these approaches have been limited to only scheduling a single-vehicle. The *multi*-extensions of these routing problems, such as *multiple* TSP and *multiple* VRP, are underrepresented in the deep learning research community, even though they capture a broader set of the real-world problems and pose a more significant scientific challenge.

The multiple traveling salesmen problem (mTSP) aims to determine a set of subroutes for each salesman, given m salesmen and N cities that need to be visited by one of the salesmen, and a depot where salesmen are initially located and to which they return. The objective of a mTSP is either minimizing the sum of subroute lengths (*MinSum*) or minimizing the length of the longest subroute

(*MinMax*). In general, the *MinMax* objective is more practical, as one seeks to visit all cities as soon as possible (i.e., total completion time minimization). In contrast, the *MinSum* formulation, in general, leads to highly imbalanced solutions where one of the salesmen visits most of the cities, which results in longer total completion time (Lupoaie et al., 2019).

In this study, we formulate (*MinMax* mTSP as a Markov Decision Process (MDP) and derive an RL-based construction heuristic that can solve mTSP instances with any numbers of salesmen and cities. Solving *MinMax* mTSP using an RL-based construction approach is challenging because one needs to train the scheduling policy that can coordinate multiple workers to complete the distributed tasks as quickly as possible using only the delayed and sparse reward signal (makespan). In addition, the state representation should be flexible and transferable to varying numbers of salesmen and cities. (coordination, transferability, sparse and delayed reward).

The proposed method first presents a state (partial solution to mTSP) as a set of graphs, each of which captures specific relationships among works, cities, and a depot. The proposed method then employs type-aware graph attention (TGA) to compute the node embeddings. Lastly, the proposed model using the node embeddings make the next assignment action. The type aware graph node embedding technique specially designed for mTSP is an essential component for deriving the cooperative and transferable scheduling policy. Additionally, to effectively train the graph representation and policy learning modules effectively with sparse and delayed reward (makespan), we propose a stable RL training scheme which significantly stabilizes the training and improves the generalization performance.

We have empirically shown that the proposed method achieves the performance comparable to Google OR-Tools, a highly optimized meta-heuristic baseline. The proposed approach outperforms OR-Tools in many cases on in-training, out-of-training problem distributions, and real-world problem instances.

2 RELATED WORK

Construction RL approaches A seminal body of work focused on the construction approach in the RL setting for solving CO problems (Bello et al., 2016; Nazari et al., 2018; Kool et al., 2018; Khalil et al., 2017). These approaches utilize *encoder-decoder* architecture, that encodes the problem structure into a hidden embedding first, and then autoregressively decodes the complete solution. Bello et al. (2016) utilized LSTM (Hochreiter & Schmidhuber, 1997) based encoder and decode the complete solution (tour) using Pointer Network (Vinyals et al., 2015) scheme. Since the routing tasks are often represented as graphs, Nazari et al. (2018) proposed an attention based encoder, while using LSTM decoder. Recently, Kool et al. (2018) proposed to use Transformer-like architecture (Vaswani et al., 2017) to solve several variants of TSP and single-vehicle CVRP. On the contrary, Khalil et al. (2017) do not use encoder-decoder architecture, but a single graph embedding model, *structure2vec* (Dai et al., 2016), that embeds a partial solution of the TSP and outputs the next city in the (sub)tour. (Kang et al., 2019) has extended *structure2vec* to random graph and employed this random graph embedding to solve identical parallel machine scheduling problems, the problem seeking to minimize the makespan by scheduling multiple machines.

Learned mTSP solvers The machine learning approaches for solving mTSP date back to Hopfield & Tank (1985). However, these approaches require per problem instance training. (Hopfield & Tank, 1985; Wacholder et al., 1989; Somhom et al., 1999). Among the recent learning methods, Kaempfer & Wolf (2018) encodes *MinSum* mTSP with a set-specialized variant of Transformer architecture that uses permutation invariant pooling layers. To obtain the feasible solution, they use a combination of the softassign method Gold & Rangarajan (1996) and a beam search. Their model is trained in a supervised setting using mTSP solutions obtained by Integer Linear Programming (ILP) solver. Hu et al. (2020) utilizes a GNN *encoder* and self-attention Vaswani et al. (2017) policy outputs a probability of assignment to each salesman per city. Once cities are assigned to specific salesmen, they use existing TSP solver, OR-Tools (Perron & Furnon), to obtain each worker’s sub-routes. Their method shows impressive scalability in terms of the number of cities, as they present results for mTSP instances with 1000 cities and ten workers. However, the trained model is not scalable in terms of the number of workers and can only solve mTSP problems with a pre-specified, fixed number of workers.

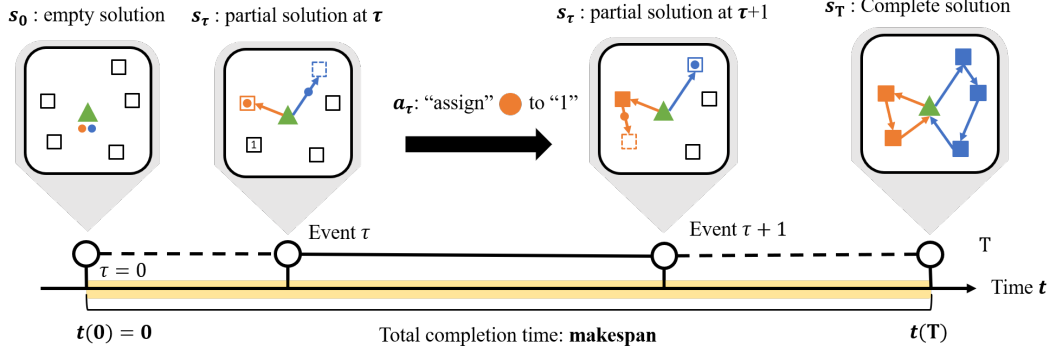


Figure 1: **The mTSP MDP** The black lined balls indicates the events of the mTSP MDP. The empty, dashed, and, filled rectangles represent the unassigned, assigned, and inactive cities, respectively. The circles represent the workers and the positions of the circles show the 2D coordinates of worker. The orange and blue colored lines shows the subtours of the orange and blue worker, respectively.

3 PROBLEM FORMULATION

We define the set of m salesmen $\mathbb{V}_T = \{1, 2, \dots, m\}$, and the set of N cities $\mathbb{V}_C = \{m+1, 2, \dots, m+N\}$. Following mTSP conventions, we define the first city as the depot. We also define the 2D-coordinates of *entities* (salesmen, cities, and the depot) as p^i . The objective of MinMax mTSP is to minimize the length of the longest subtour of salesmen, while subtours covers all cities and all subtours of salesmen end at the depot. For the clarity of explanation, we will refer to salesman as a *workers*, and cities as a *tasks*.

3.1 MDP FORMULATION FOR MINMAX MTSP

In this paper, the objective is to construct an optimal solution with a construction RL approach. Thus, we cast the solution construction process of *MinMax* mTSP as a Markov decision process (MDP). The components of the proposed MDP are as follows.

Transition. The proposed MDP transits based on *events*. We define an event as the case where any worker reaches its assigned city. We enumerate the event with the index τ for avoiding confusion from the elapsed time of the mTSP problem. $t(\tau)$ is a function that returns the time of event τ . In the proposed event-based transition setup, the state transitions coincide with the sequential expansion of the partial scheduling solution.

State. Each entity i has its own state $s_\tau^i = (p_\tau^i, \mathbb{1}_\tau^{\text{active}}, \mathbb{1}_\tau^{\text{assigned}})$ at the τ -th event. the coordinates p_τ^i is time-dependent for workers and static for tasks and the depot. Indicator $\mathbb{1}_\tau^{\text{active}}$ describes whether the entity is *active* or *inactive*. In case of tasks, inactive indicates that the task is already visited; in case of worker, inactive means that worker returned to the depot. Similarly, $\mathbb{1}_\tau^{\text{assigned}}$ indicates whether worker is assigned to a task or not. We also define the environment state s_τ^{env} that contains the current time of the environment, and the sequence of tasks visited by each worker, i.e., partial solution of the mTSP. The state s_τ of the MDP at the τ -th event becomes $s_\tau = (\{s_\tau^i\}_{i=1}^{m+N}, s_\tau^{\text{env}})$. The first state s_0 corresponds to the empty solution of the given problem instance, i.e., no cities have been visited, and all salesmen are in the depot. The terminal state s_T corresponds to a *complete* solution of the given mTSP instance, i.e., when every task has been visited, and every worker returned to the depot (See Figure 1).

Action. A scheduling action a_τ is defined as the *worker-to-task assignment*, i.e. salesman has to visit the assigned city.

Reward. We formulate the problem in a delayed reward setting. Specifically, the sparse reward function is defined as $r(s_\tau) = 0$ for all non-terminal events, and $r(s_T) = t(T)$, where T is the

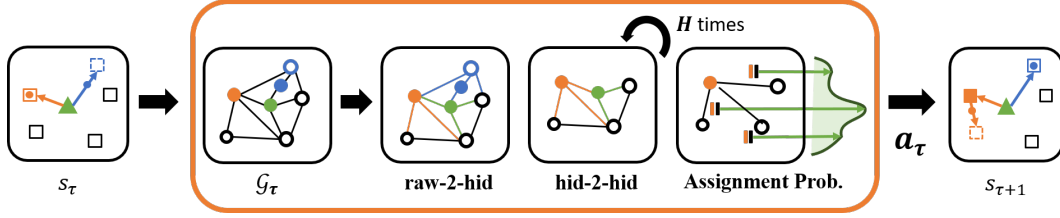


Figure 2: Assignment action determination step of ScheduleNet

index of the terminal state. In other words, a single reward signal, which is obtained only for the terminal state, is equals to the makespan of the problem instance.

4 SCHEDULENET

Given the MDP formulation for *MinMax* mTSP, we propose ScheduleNet that can recommend a scheduling action a_τ given the current state \mathcal{G}_τ represented as a graph, i.e., $\pi_\theta(a_\tau|\mathcal{G}_\tau)$. The ScheduleNet first presents a state (partial solution of mTSP) as a set of graphs, each of which captures specific relationships among workers, tasks, and a depot. Then ScheduleNet employs type-aware graph attention (TGA) to compute the node embeddings and use the computed node embeddings to determine the next assignment action (See figure 2).

4.1 WORKER-TASK GRAPH REPRESENTATION

Whenever an event occurs and the global state s_τ of the MDP is updated at τ , ScheduleNet constructs a directed complete graph $\mathcal{G}_\tau = (\mathbb{V}, \mathbb{E})$ out of s_τ , where $\mathbb{V} = \mathbb{V}_T \cup \mathbb{V}_C$ is the set of nodes and \mathbb{E} is the set of edges. We drop the time iterator τ to simplify the notations since the following operations only for the given time step. The nodes and edges and their associated features are defined as:

- v_i denotes the node corresponding entity i in mTSP problem. The node feature x_i for v_i is equal to the state s_τ^i of entity i . In addition, k_i denote the type of node v_i . For instance, if the entity i is *worker* and its $\mathbb{1}_\tau^{active} = 1$, then the k_i becomes *active-worker* type.
- e_{ij} denotes the edge between between source node v_i and destination node v_j , representing the relationships between the two. The edge feature w_{ij} is equal to the Euclidean distance between the two nodes.

4.2 TYPE-AWARE GRAPH ATTENTION EMBEDDING

In this section, we describe a *type-aware* graph attention (TGA) embedding procedure. We denote h_i and h_{ij} as the node and the edge embedding, respectively, at a given time step, and h'_i and h'_{ij} as the updated embedding by TGA embedding. A single iteration of TGA embedding consists of three phases: (1) edge update, (2) message aggregation, and (3) node update.

Type-aware Edge update Given the node embeddings h_i for $v_i \in \mathbb{V}$ and the edge embeddings h_{ij} for $e_{ij} \in \mathbb{E}$, ScheduleNet computes the updated edge embedding h'_{ij} and the attention logit z_{ij} as:

$$\begin{aligned} h'_{ij} &= \text{TGA}_{\mathbb{E}}([h_i, h_j, h_{ij}], k_j) \\ z_{ij} &= \text{TGA}_{\mathbb{A}}([h_i, h_j, h_{ij}], k_j) \end{aligned} \quad (1)$$

where $\text{TGA}_{\mathbb{E}}$ and $\text{TGA}_{\mathbb{A}}$ are, respectively, the type aware edge update function and the type aware attention function, which are defined for the specific type k_j of the source node v_j . The updated edge feature h'_{ij} can be thought of as the message from the source node v_j to the destination node v_i , and the attention logit z_{ij} will be used to compute the importance of this message.

In computing the updated edge feature (message), $\text{TGA}_{\mathbb{E}}$ and $\text{TGA}_{\mathbb{A}}$ first compute the “type-aware” edge encoding u_{ij} , which can be seen as a dynamic edge feature varying depending on the source node type, to effectively model the complex type-aware relationships among the nodes. Using the

computed “type-aware” edge encoding u_{ij} , these two functions then compute the updated edge feature and attention logit using a multiplicative interaction (MI) layer (Jayakumar et al., 2019). The use of MI layer significantly reduces the number of parameters to learn without discarding the expressibility of the embedding procedure. The detailed architecture for $\text{TGA}_{\mathbb{E}}$ and $\text{TGA}_{\mathbb{A}}$ are provided in Appendix A.1.

Type-aware Message aggregation The distribution of the node types in the mTSP graphs is highly imbalanced, i.e., the number of task-specific node types is much larger than the worker specific ones. This imbalance is problematic, specifically, during the message aggregation of GNN, since permutation invariant aggregation functions are akin to ignore messages from few-but-important nodes in the graph. To alleviate such an issue, we propose the following *type-aware* message aggregation scheme.

We first define the type k neighborhood of node v_i as the set of the k -type source nodes that are connected to the destination node v_i , i.e., $\mathcal{N}_k(i) = \{v_j | k_j = k, \forall l \in \mathcal{N}(i)\}$, where $\mathcal{N}(i)$ is the neighborhood set of node v_i containing the nodes that are connected to node v_i with edges.

The node v_i aggregates separately messages from the same type of source nodes. For example, the aggregated message m_i^k from k -type source nodes is computed as:

$$m_i^k = \sum_{j \in \mathcal{N}_k(i)} \alpha_{ij} h'_{ij} \quad (2)$$

where α_{ij} is the attention score computed using the attention logits computed before as:

$$\alpha_{ij} = \frac{\exp(z_{ij})}{\sum_{j \in \mathcal{N}_k(i)} \exp(z_{ij})} \quad (3)$$

Finally, all aggregated messages per type are concatenated to produce the total aggregated message m_i for node v_i as

$$m_i = \text{concat}(\{m_i^k | k \in \mathbb{K}\}) \quad (4)$$

Type-aware Node update The aggregated message m_i for node v_i is then used to compute the updated node embedding h'_i using the type-aware graph node update function $\text{TGA}_{\mathbb{V}}$ as:

$$h'_i = \text{TGA}_{\mathbb{V}}(h_i, m_i, k_i) \quad (5)$$

Being different from $\text{TGA}_{\mathbb{E}}$ and $\text{TGA}_{\mathbb{A}}$, $\text{TGA}_{\mathbb{V}}$ is defined per different destination node type k_i .

4.3 ASSIGNMENT PROBABILITY COMPUTATION

ScheduleNet model consists of two type-aware graph embedding layers that utilize the embedding procedure explained in the section above. The first embedding layer *raw-2-hid* is used to encode initial node and edge features x_i and w_{ij} of the (full) graph \mathcal{G}_τ , to obtain initial hidden node and edge features $h_i^{(0)}$ and $h_{ij}^{(0)}$, respectively.

We define the *target subgraph* \mathcal{G}_τ^s as the subset of nodes and edges from the original (full) graph \mathcal{G}_τ that only includes a target-worker (unassigned-worker) node and all unassigned-city nodes. The second embedding layer *hid-2-hid* embeds the *target* subgraph \mathcal{G}_τ^s H times. In other words, a hidden node and edge embeddings $h_i^{(0)}$ and $h_{ij}^{(0)}$ are iteratively updated H times to obtain final hidden embeddings $h_i^{(H)}$ and $h_{ij}^{(H)}$, respectively. The final hidden embeddings are then used to make decision regarding the *worker-to-task assignment*.

Specifically, probability of assigning target worker i to task j is computed as

$$\begin{aligned} y_{ij} &= \text{MLP}_{\text{actor}}([h_i^{(H)}; h_j^{(H)}; h_{ij}^{(H)}]) \\ p_{ij} &= \text{softmax}(\{y_{ij}\}_{j \in \mathbb{A}(\mathcal{G}_\tau)}) \end{aligned} \quad (6)$$

where the $h_i^{(H)}$, and $h_{ij}^{(H)}$ is the final hidden node, edge embeddings, respectively. In addition, $\mathbb{A}(\mathcal{G}_\tau)$ denote the set of feasible actions defined as $\{v_j | k_j = \text{“Unassigned-task”} \forall j \in \mathbb{V}\}$.

5 TRAINING SCHEDULENET

In this section, we describe the training scheme of the ScheduleNet. Firstly, we explain reward normalization scheme which is used to reduce the variance of the reward. Secondly, we introduce a stable RL training scheme which significantly stabilizes the training process.

Makespan normalization

As mentioned in Section 3.1, we use the makespan of mTSP as the only reward signal for training RL agent. We denote the makespan of given policy π as $M(\pi)$. We observe that, the makespan $M(\pi)$ is a highly volatile depending on the problem size (number of cities and salesmen), the topology of the map, and the policy. To reduce the variance of the reward, we propose the following normalization scheme:

$$m(\pi, \pi_b) = \frac{M(\pi_b) - M(\pi)}{M(\pi_b)} \quad (7)$$

where π and π_b is the evaluation and baseline policy, respectively.

The normalized makespan $m(\pi, \pi_b)$ is similar to (Kool et al., 2018), but we additionally divide the performance difference by the makespan of the baseline policy, which further reduces the variance that is induced by the size of the mTSP instance.

From the normalized terminal reward $m(\pi, \pi_b)$, we compute the normalized return as follows:

$$G_\tau(\pi, \pi_b) := \gamma^{T-\tau} m(\pi, \pi_b) \quad (8)$$

where T is the index of the terminal state, and γ is the discount factor.

Stable RL training It is well known that the solution quality of CO problems, including the makespan of mTSP, is extremely sensitive to the action selection, and it thus prevents the stable policy learning. To address this problem, we propose the *clipped* REINFORCE, a variant PPO *without* the learned value function. We empirically found that it is hard to train the value function¹, thus we use normalized returns $G_\tau(\pi_\theta, \pi_b)$ directly. Then, the objective of the clipped REINFORCE is given as follows:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathcal{G}_\tau, a_\tau) \sim \pi_\theta} [\min(\text{clip}(\rho_\tau, 1 - \epsilon, 1 + \epsilon) G_\tau(\pi_\theta, \pi_b), \rho_\tau G_\tau(\pi_\theta, \pi_b))] \quad (9)$$

where

$$\rho_\tau = \frac{\pi_\theta(a_\tau | \mathcal{G}_\tau)}{\pi_{\theta_{old}}(a_\tau | \mathcal{G}_\tau)} \quad (10)$$

and $(\mathcal{G}_\tau, a_\tau) \sim \pi_\theta$ is the state-action marginal following π_θ , and $\pi_{\theta_{old}}$ is the old policy.

After updating the policy π_θ , we smooth the parameters of policy π_θ with the Polyak average (Polyak & Juditsky, 1992) to further stabilize policy training. The training procedure of ScheduleNet is given in Appendix A.2.1.

6 EXPERIMENTS

We train the ScheduleNet using mTSP instances whose number m of tasks and the number N of workers are sampled from $m \sim U(2, 4)$ and $N \sim U(10, 20)$, respectively. This trained ScheduleNet policy is then evaluated on the various dataset, including randomly generated uniform mTSP datasets, mTSPLib, and randomly generated uniform TSP dataset, TSPLib, and TSP (dai). Unless explicitly mentioned, all performance results are obtained from this single trained model. See Appendix for further training details.

¹Note that the value function is trained to predict the makespan of the state to serve as an advantage estimator. Due to the combinatorial nature of the mTSP, the target of value function, makespan, is highly volatile, which makes training value function hard.

Table 1: MTSP Uniform Results.

	$N = 20$			$N = 50$		
	$m = 2$	$m = 3$	$m = 5$	$m = 5$	$m = 7$	$m = 10$
Obj. (OR-Tools)	2.527	2.083	1.833	2.076	1.946	1.937
Obj. (ScheduleNet)	2.744	2.260	1.959	2.352	2.172	2.098
Obj. (Hu et. al.)	–	–	–	2.115	–	1.965
Gap (vs OR-Tools)	1.087	1.086	1.071	1.135	1.122	1.086
	$N = 100$			$N = 200$		
	$m = 5$	$m = 10$	$m = 15$	$m = 10$	$m = 15$	$m = 20$
Obj. (OR-Tools)	2.452	1.993	2.031	?	2.295	2.370
Obj. (ScheduleNet)	2.892	2.251	2.185	?	2.248	2.299
Obj. (Hu et. al.)	2.480	2.087	–	–	–	–
Gap (vs OR-Tools)	1.179	1.129	1.075	1.164(?)	1.011	0.971

Table 2: MTSPLIB CITE MINMAX Results Result with \dagger 1, 2, 3, 4 were obtained by running CPLEX for 120, 96, 168, 216 hours, respectively. The optimal result(s) are indicated with *, otherwise, CPLEX results are reported as the average of the upper and lower bound.

Instance	m	CPLEX	OR-Tools	ScheduleNet	SOM	ACO	EA
<i>eil51</i>	2	222.73*	243.02	259.67	278.44	248.76	276.62
	3	155.13	170.05	172.16	210.25	180.59	208.16
	5	110.43	127.50	118.94	157.68	135.09	151.21
	7	92.44 $^{\dagger 1}$	112.07	112.42	136.84	119.96	123.88
<i>berlin52</i>	2	4079.63	4665.47	4816.30	5350.83	4388.99	5038.33
	3	2999.01	3311.31	3372.14	4197.61	3468.9	3865.45
	5	2056.54	2482.57	2615.57	3461.93	2733.56	2853.63
	7	1856.49	2440.92	2576.04	3125.21	2510.09	2543.73
<i>eil76</i>	2	280.85	318.00	334.10	364.02	308.53	365.72
	3	191.84	212.41	226.54	278.63	224.56	285.43
	5	133.95	143.38	168.03	210.69	163.93	211.91
	7	113.985	128.31	151.31	183.09	146.88	177.83
<i>rat99</i>	2	674.85	762.19	789.98	927.36	767.15	896.72
	3	512.13 $^{\dagger 2}$	552.09	579.28	756.08	620.45	739.43
	5	402.71 $^{\dagger 3}$	473.66	502.49	624.38	525.54	596.87
	7	373.80 $^{\dagger 4}$	442.47	471.67	564.14	492.13	534.91
Gap		1.00	1.145	1.206	1.464	1.216	1.287

6.1 MTSP RESULTS

Random mTSP results We firstly investigate the generalization performance of ScheduleNet on the randomly generated uniform maps with varying numbers of tasks and workers. We report the results of OR-Tools, as well as (Hu et al., 2020) when available. Table 1 shows that ScheduleNet in overall produces a slightly longer makespan than OR-Tools even for the large-sized mTSP instances. As the complexity of the target mTSP instance increases, the gap between ScheduleNet and OR-Tools decreases, even showing the cases where ScheduleNet outperforms OR-Tools. The result empirically proves that ScheduleNet, even trained with small-sized mTSP instances, can solve large scale problems well (generalization and scalability).

mTSPLib results The trained ScheduleNet is employed to solve the benchmark problems in mT-SPLib (cite), without parameter retraining, to validate the generalization capability of ScheduleNet on unseen mTSP instances whose problem structures can be completely different from the ones used during training. Table 2 compares the performance of the ScheduleNet to other baseline models, including CPLEX (optimal solution), OR-Tools, and other heuristics. In general, ScheduleNet

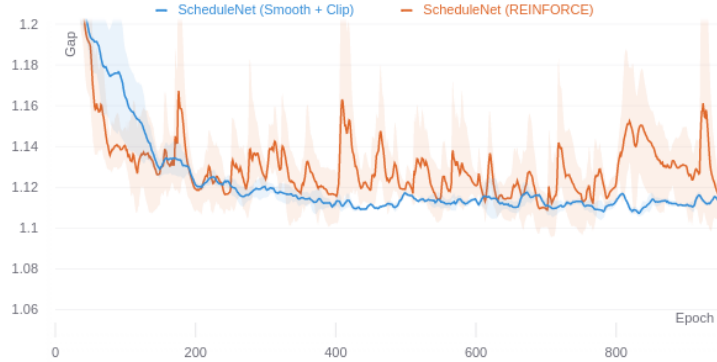


Figure 3: **Training curves** The orange and blue line shows the averaged training performance of ScheduleNet trained with REINFORCE (Sutton et al., 2000) and clipped REINFORCE with parameter smoothing, respectively. The shadow regions indicate the standard deviations of the models. We repeat ten times with different random seeds for each setup.

produces a slightly longer makespan than OR-Tools but significantly lower makespan than other heuristics. It is worth noting that Our model is the first one that has empirically proven that an RL-based contraction heuristic generalizes well over unseen mTSP whose problem distribution is entirely different from the training instances.

6.2 EFFECTIVENESS OF THE PROPOSED TRAINING SCHEME

Figure 3 compares the training curves of ScheduleNet model that was trained with REINFORCE-only against the performance of the ScheduleNet with clipped REINFORCE with parameter smoothing. This clearly illustrates how the parameter smoothing and policy clipping results in a considerably more stable training.

7 CONCLUSION

We proposed ScheduleNet for solving *MinMax* mTSP, the problem seeking to minimize the total completion time for multiple workers to complete the geographically distributed tasks. The use of type-aware graphs and the specially designed TGA graph node embedding allows the trained ScheduleNet policy to induce the coordinated strategic subroutes of the workers and to be well transferred to unseen mTSP with any numbers of workers and tasks. We have empirically shown that the proposed method achieves the performance comparable to Google OR-Tools, a highly optimized meta-heuristic baseline. All in all, this study has shown the potential that the proposed ScheduleNet can be effectively used to schedule multiple vehicles for solving large-scale, practical, real-world applications.

REFERENCES

- Sungsoo Ahn, Younggyo Seo, and Jinwoo Shin. Deep auto-deferring policy for combinatorial optimization. 2019.
- Irwan Bello, Hieu Pham, Quoc V Le, Mohammad Norouzi, and Samy Bengio. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.
- Yoshua Bengio, Andrea Lodi, and Antoine Prouvost. Machine learning for combinatorial optimization: a methodological tour d’horizon. *European Journal of Operational Research*, 2020.
- Xinyun Chen and Yuandong Tian. Learning to perform local rewriting for combinatorial optimization. In *Advances in Neural Information Processing Systems*, pp. 6281–6292, 2019.
- Hanjun Dai, Bo Dai, and Le Song. Discriminative embeddings of latent variable models for structured data. In *International conference on machine learning*, pp. 2702–2711, 2016.

- Steven Gold and Anand Rangarajan. Softassign versus softmax: Benchmarks in combinatorial optimization. In *Advances in neural information processing systems*, pp. 626–632, 1996.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- John J Hopfield and David W Tank. “neural” computation of decisions in optimization problems. *Biological cybernetics*, 52(3):141–152, 1985.
- Yujiao Hu, Yuan Yao, and Wee Sun Lee. A reinforcement learning approach for optimizing multiple traveling salesman problems over graphs. *Knowledge-Based Systems*, 204:106244, 2020.
- Siddhant M Jayakumar, Wojciech M Czarnecki, Jacob Menick, Jonathan Schwarz, Jack Rae, Simon Osindero, Yee Whye Teh, Tim Harley, and Razvan Pascanu. Multiplicative interactions and where to find them. In *International Conference on Learning Representations*, 2019.
- Yoav Kaempfer and Lior Wolf. Learning the multiple traveling salesmen problem with permutation invariant pooling networks. *arXiv preprint arXiv:1803.09621*, 2018.
- Hyunwook Kang, Aydar Mynbay, James Morrison, and Jinkyoo Park. Learning scalable and transferable multi-robot/machine sequential assignment planning via graph embedding. *arXiv preprint arXiv:1905.12204*, 2019.
- Elias Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, and Le Song. Learning combinatorial optimization algorithms over graphs. In *Advances in Neural Information Processing Systems*, pp. 6348–6358, 2017.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Wouter Kool, Herke Van Hoof, and Max Welling. Attention, learn to solve routing problems! *arXiv preprint arXiv:1803.08475*, 2018.
- Vlad-Ioan Lupoae, Ivona-Alexandra Chili, Mihaela Elena Breaban, and Madalina Raschip. Som-guided evolutionary search for solving minmax multiple-tsp. In *2019 IEEE Congress on Evolutionary Computation (CEC)*, pp. 73–80. IEEE, 2019.
- Nina Mazyavkina, Sergey Sviridov, Sergei Ivanov, and Evgeny Burnaev. Reinforcement learning for combinatorial optimization: A survey. *arXiv preprint arXiv:2003.03600*, 2020.
- Mohammadreza Nazari, Afshin Oroojlooy, Lawrence Snyder, and Martin Takác. Reinforcement learning for solving the vehicle routing problem. In *Advances in Neural Information Processing Systems*, pp. 9839–9849, 2018.
- Laurent Perron and Vincent Furnon. Or-tools. URL <https://developers.google.com/optimization/>.
- Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992.
- Samerkae Somhom, Abdolhamid Modares, and Takao Enkawa. Competition-based neural network for the multiple travelling salesmen problem with minmax objective. *Computers & Operations Research*, 26(4):395–407, 1999.
- Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pp. 1057–1063, 2000.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Advances in neural information processing systems*, pp. 2692–2700, 2015.
- E Wacholder, J Han, and RC Mann. A neural network algorithm for the multiple traveling salesmen problem. *Biological Cybernetics*, 61(1):11–19, 1989.

A APPENDIX

A.1 DETAILS OF TYPE-AWARE GRAPH ATTENTION EMBEDDING

In this section, we thoroughly describe a *type-aware* graph embedding procedure. Similar to the main body, We overload notations for the simplicity of notation such that the input node and edge feature as h_i and h_{ij} , and the embedded node and edge feature h'_i and h'_{ij} , respectively.

The proposed graph embedding step consists of three phases: (1) type-aware edge update, (2) type-aware message aggregation, and (3) type-aware node update.

Type-aware Edge update The edge update scheme is designed to reflect the complex type relationship among the entities while updating edge features. First the *context* embedding c_{ij} of edge e_{ij} computed using the source node type, k_i , such that:

$$c_{ij} = \text{MLP}_{etype}(k_i) \quad (11)$$

where MLP_{etype} is the edge type encoder. The source node types are embedded into the context embedding c_{ij} using MLP_{etype} . Next, the type-aware edge encoding u_{ij} is computed using the Multiplicative Interaction (MI) layer (Jayakumar et al., 2019) as follows:

$$u_{ij} = \text{MI}_{edge}([h_i; h_j; h_{ij}], c_{ij}) \quad (12)$$

where MI_{edge} is the edge MI layer. We utilize the MI layer, which dynamically generates its parameter depending on the context c_{ij} and produces and produces “type-aware” edge encoding u_{ij} , to effectively model the complex type relationships among the nodes. “type-aware” edge encoding u_{ij} can be seen as a dynamic edge feature which varies depending on the source node type. After the updated edge embedding h'_{ij} and its attention logit z_{ij} is obtained as:

$$h'_{ij} = \text{MLP}_{edge}(u_{ij}) \quad (13)$$

$$z_{ij} = \text{MLP}_{attn}(u_{ij}) \quad (14)$$

where MLP_{edge} and MLP_{attn} is the edge updater and logit function, respectively. the edge updater and logit function produces updated edge embedding and logits from the “type-aware” edge.

The computation steps of equation 11, 12, and ?? are defined as $\text{TGA}_{\mathbb{B}}$. Similarly, the computation steps of equation 11, 12, and ?? are defined as $\text{TGA}_{\mathbb{A}}$.

Message aggregation First, we define the type- k neighborhood of node v_i such that $\mathcal{N}_k(i) = \{e_{li} | k_l = k, \forall l \in \mathcal{N}(i)\}$, where $\mathcal{N}(i)$ is the neighborhood set of node i . The type- k neighborhood is the set of edges heading to node i , and their source nodes have type k . The proposed type-aware message aggregation procedure computes attention score α_{ji} for the e_{ji} , which starts from node j and heads to node i , such that:

$$\alpha_{ji} = \frac{\exp(z_{ji})}{\sum_{l \in \mathcal{N}_{k_j}(i)} \exp(z_{li})} \quad (15)$$

Intuitively speaking, The proposed attention scheme normalizes the attention logits of incoming edges over the types. Therefore, the attention scores sum up to 1 over each type- k neighborhood. Next, the type- k neighborhood message $m_{i,k}$ for node v_i is computed as:

$$m_{i,k} = \sum_{l \in \mathcal{N}_k(i)} \alpha_{li} h'_{li} \quad (16)$$

In this aggregation step, the incoming messages of node i are aggregated type-wisely. Finally, all incoming type neighborhood messages are concatenated to produce (inter-type) aggregated message m_i for node v_i , such that:

$$m_i = \text{concat}(\{m_{i,k} | k \in \mathbb{K}\}) \quad (17)$$

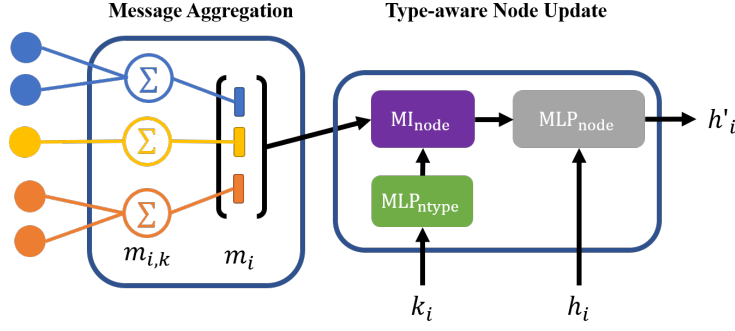


Figure 4: **Type-aware graph attention embedding** We omit the type-aware edge update for the clarity of visualization.

Node update Similar to the edge update phase, first, the context embedding c_i is computed for each node v_i :

$$c_i = \text{MLP}_{ntype}(k_i) \quad (18)$$

where MLP_{ntype} is the node type encoder. Then, the updated hidden node embedding h'_i is computed as below:

$$h'_i = \text{MLP}_{node}([h_i; u_i]) \quad (19)$$

where $u_i = \text{MI}_{node}(m_i, c_i)$ is the type-aware node embedding that is produced by MI_{node} layer using aggregated messages m_i and the context embedding c_i .

The computation steps of equation 18, and 19 are defined as $\text{TGA}_{\mathbb{E}}$. The overall computation procedure TGA is illustrated in Figure 4.

A.2 DETAILS OF SCHEDULENET TRAINING

A.2.1 TRAINING PSEUDO CODE

In this section, we presents a pseudocode for training ScheduleNet.

Algorithm 1: ScheduleNet Training

Input: Training policy π_θ

Output: Smoothed policy π_ϕ

```

1 Initialize smoothed policy with parameters  $\phi \leftarrow \theta$ .
2 for update step do
3   Generate a random mTSP instance  $I$ 
4   for number of episodes do
5     Construct mTSP MDP from the instance  $I$ 
6      $\pi_b \leftarrow \arg \max(\pi_\theta)$ 
7     Collect samples with  $\pi_\theta$  and  $\pi_b$  from the mTSP MDP.
8    $\pi_{\theta \text{old}} \leftarrow \pi_\theta$ 
9   for inner updates  $K$  do
10     $\theta \leftarrow \theta + \alpha \nabla_\theta \mathcal{L}(\theta)$ 
11   $\phi \leftarrow \beta \phi + (1 - \beta) \theta$ 

```

A.2.2 HYPERPARAMETERS

In this section, we fully explain hyperparameters of SchduleNet. **Network Architecture** We use the same hyperparameters for the *raw-2-hid* TGA layer and the *hid-2-hid* TGA layer. MLP_{etype} and MLP_{ntype} has one hidden layer with 32 neurons and their output dimensions are both 32. Both MI layers has 64 dimensional outputs. MLP_{edge} , MLP_{attn} , and MLP_{node} has 2 hidden layers with 32

neurons. MLP_{actor} has 2 hidden layers and the hidden layers has 128 neurons each. We use ReLU activation functions for all hidden layers. The hidden graph embedding step H is two.

Training We use Adam (Kingma & Ba, 2014) with learning rate value of 0.001. Inner updates steps K is three. The smoothing parameter β is 0.95.

A.3 TRANSFERABILITY TEST ON TSP ($m = 1$)

The trained ScheduleNet has been employed to solve random TSP instances. Because ScheduleNet can be used to schedule any m number of workers, if we set $m = 1$, it can be used to schedule TSP instance without further training. Table 3 shows the results on this transferability experiments.

Table 3 shows that the trained ScheduleNet can solve reasonably well on random TSP instances, although ScheduleNet has never been exposed to such TSP instances. Note that as the size of TSP increases, the gap between the ScheduleNet and other models becomes smaller. If the ScheduleNet is trained with TSP instances with $m = 1$, the performance can be further improved. However, we did not try that experiment to check its transferability over different types of routing problems with different objectives.

Table 3: Performance comparison on random uniform TSP instances. The Obj. defines the makespan, i.e., the length of the tour of the salesman.

Method	$N = 20$			$N = 50$			$N = 100$		
	Obj.	Gap	Time	Obj.	Gap	Time	Obj.	Gap	Time
Concorde	3.84	0.00%	(1m)	5.70	0.00%	(2m)	7.76	0.00%	(3m)
LKH3	3.84	0.00%	(18s)	5.70	0.00%	(5m)	7.76	0.00%	(21m)
Gurobi	3.84	0.00%	(7s)	5.70	0.00%	(2m)	7.76	0.00%	(17m)
Nearest Insertion	4.33	12.91%	(1s)	6.78	19.03%	(2s)	9.46	21.82%	(6s)
Farthest Insertion	3.93	2.36%	(1s)	6.01	5.53%	(2s)	8.35	7.59%	(7s)
OR-Tools	3.85	0.37%	(1s)	5.80	1.83%	(2s)	7.99	2.90%	(7s)
Bello et al. (2016)	3.89	1.42%	–	5.95	4.46%	–	8.30	6.90%	–
Khalil et al. (2017)	3.89	1.42%	–	5.99	5.16%	–	8.31	7.03%	–
AM (greedy)	3.85	0.34%	(0s)	5.80	1.76%	(2s)	8.12	4.53%	(6s)
Am (sampling)	3.84	0.08%	(5m)	5.73	0.52%	(24m)	7.94	2.26%	(1h)
ScheduleNet	4.01	4.34%	(2s)	6.11	7.01%	(9s)	8.71	12.35%	(34s)