
RL4CO: an Extensive Reinforcement Learning for Combinatorial Optimization Benchmark

Federico Berto^{*1}, Chuanbo Hua^{*1}, Junyoung Park^{*1}, Laurin Luttmann^{*2},
Yining Ma³, Fanchen Bu¹, Jiarui Wang⁴, Haoran Ye⁵, Minsu Kim¹,
Sanghyeok Choi¹, Nayeli Gast Zepeda⁶, André Hottung⁶, Jianan Zhou³, Jieyi Bi³,
Yu Hu⁷, Fei Liu⁸, Hyeonah Kim¹, Jiwoo Son¹⁵, Haeyeon Kim¹,
Davide Angioni⁹, Wouter Kool¹⁰, Zhiguang Cao¹¹, Qingfu Zhang⁸, Joungho Kim¹,
Jie Zhang³, Kijung Shin¹, Cathy Wu¹², Sungsoo Ahn¹³, Guojie Song⁵
Changhyun Kwon^{1,14}, Kevin Tierney⁶, Lin Xie^{2,14}, Jinkyoo Park^{1,15}

¹KAIST, ²Leuphana University, ³Nanyang Technological University,
⁴Southeastern University, ⁵Peking University, ⁶Bielefeld University,
⁷Soochow University, ⁸City University of Hong Kong, ⁹University of Brescia,
¹⁰ORTEC, ¹¹Singapore Management University, ¹²MIT, ¹³POSTECH,
¹⁴Twente University, ¹⁵OMELET, AI4CO[🔥]

Abstract

1 Deep reinforcement learning (RL) has recently shown significant benefits in solv-
2 ing combinatorial optimization (CO) problems, reducing reliance on domain ex-
3 pertise, and improving computational efficiency. However, the field lacks a uni-
4 fied benchmark for easy development and standardized comparison of algorithms
5 across diverse CO problems. To fill this gap, we introduce RL4CO, a unified
6 and extensive benchmark with in-depth library coverage of 23 state-of-the-art
7 methods and more than 20 CO problems. Built on efficient software libraries
8 and best practices in implementation, RL4CO features modularized implemen-
9 tation and flexible configuration of diverse RL algorithms, neural network archi-
10 tectures, inference techniques, and environments. RL4CO allows researchers to
11 seamlessly navigate existing successes and develop their unique designs, facili-
12 tating the entire research process by decoupling science from heavy engineering.
13 We also provide extensive benchmark studies to inspire new insights and future
14 work. RL4CO has attracted numerous researchers in the community and is open-
15 sourced at <https://github.com/ai4co/rl4co>.

16 1 Introduction

17 Combinatorial optimization (CO) focuses on finding optimal solutions for problems with discrete
18 variables and has broad applications, including vehicle routing [89, 60], scheduling [128], and hard-
19 ware device placement [53]. Given that the combinatorial space expands exponentially and exhibits
20 NP-hard characteristics, the operations research (OR) community has traditionally tackled these
21 challenges through the development of mathematical programming algorithms [35] and handcrafted
22 heuristics [27]. Despite their success, these methods still face significant limitations: mathemat-
23 ical programming struggles with scaling, while handcrafted heuristics require significant domain-
24 specific adjustments for different CO problems.

*Equal contribution.

🔥 Work made with contributions from the AI4CO open research community.

25 Recently, to address these limitations, neural combinatorial optimization (NCO) [7] has emerged.
26 It employs deep neural networks to automate the problem-solving process and significantly reduces
27 the computation demands and the need for domain expertise. Recent NCO works mainly leverage
28 the reinforcement learning (RL) paradigm, making significant strides in improving exploration ef-
29 ficiency [62, 54], relaxing the needs of obtaining optimal solutions, and extending to various CO
30 tasks [128, 89, 60, 53]. Although supervised learning (SL) methods [29] are shown to be effective in
31 NCO, they require the availability of high-quality solutions, which is unrealistic for large instances
32 or theoretically hard problems. Therefore, we focus on the widespread RL paradigm in this paper.

33 Despite the growing popularity and advancements in using reinforcement learning for solving com-
34 binatorial optimization, there remains a lack of a unified benchmark for analyzing past works under
35 consistent implementations and conditions. The absence of a standardized benchmark hinders NCO
36 researchers’ efforts to make impactful advancements and leverage existing successes, as it becomes
37 challenging to determine the superiority of one method over another. Moreover, the significance of
38 NCO lies in its potential for generalizability across multiple problems without extensive problem-
39 specific knowledge. Variations in implementation can make it difficult for new researchers to en-
40 gage with the NCO community, and inconsistent comparisons obstruct straightforward performance
41 evaluations. These issues pose significant challenges and underscore the need for a comprehensive
42 benchmark to streamline research and foster consistent progress.

43 **Contributions.** To bridge this gap, we introduce RL4CO, the first comprehensive benchmark with
44 multiple baselines, environments, and boilerplate from the literature, all implemented in a *modular*,
45 *flexible*, *accelerated*, and *unified* manner. Our aim is to facilitate the entire research process for
46 the NCO community with the following key contributions: 1) **Simplifying development** through
47 modularizing 27 environments and 23 existing baseline models, allowing for flexible and automated
48 combinations for effortless testing, switching, and achieving state-of-the-art performance; 2) **En-**
49 **hancing the training and testing efficiency** through the customized unified pipeline tailored for
50 the NCO community based on advanced libraries such as TorchRL [15], PyTorch Lightning [31],
51 Hydra [123], and TensorDict [15]; 3) **Standardizing evaluation** to ensure fair and comprehensive
52 comparisons, enabling researchers to automatically test a broader range of problems from diverse
53 distributions and gather valuable insights using our testbed. Overall, RL4CO eliminates the need
54 for repetitive heavy engineering in the NCO community and fosters seamless future development by
55 building on existing successes, enabling advanced innovation and progress in the field.

56 2 Related Works

57 **Neural Combinatorial Optimization.** Neural combinatorial optimization (NCO) utilizes machine
58 learning techniques to automatically develop novel heuristics for solving NP-hard CO problems.
59 We classify the majority of NCO research from the following perspectives: 1) *Learning Paradigms*:
60 researchers have employed supervised learning [115, 108, 29, 75] to approximate optimal solutions
61 to CO instances. Further research leverages reinforcement learning [6, 89, 60, 62], and unsupervised
62 learning [39, 84] to ease the difficulty of obtaining (near-)optimal solutions. 2) *Models*: various deep
63 learning architectures such as recurrent neural networks [115, 22, 68], graph neural networks [48,
64 84], Transformers [60, 62], diffusion models [108], and GFlowNets [129, 56] have been employed.
65 3) *Problems*: NCO has demonstrated great success in various problems, including vehicle routing
66 problems (VRPs) (e.g., traveling salesman problem and capacitated VRP), scheduling problems
67 (e.g., job shop scheduling problems [128]), hardware device placement [53], and graph-based CO
68 problems (e.g., maximum independent set [23, 2] and maximum cut [129]). 4) *Heuristic Types*:
69 generally, the learned heuristics can be categorized as *constructive* in an autoregressive [60] or non-
70 autoregressive [48] way, and *improvement* heuristics, which leverage traditional heuristics [120, 80]
71 and meta-heuristics [105]. We refer to Bengio et al. [7] for a comprehensive survey. In this paper,
72 we focus on the reinforcement learning paradigm due to its effectiveness and flexibility. Notably,
73 the proposed RL4CO is versatile to support most combinations of models, problems and heuristic
74 types, making it an apt library and benchmark for future research in NCO.

Table 1: Comparison of libraries in reinforcement learning for combinatorial optimization.

Library	Environments #	Baselines [†] #	Hardware Acceleration	Availability	Modular Baselines	Open Community
ORL [4]	3	1	×	×	×	×
OR-Gym [42]	9	-	×	✓	×	×
Graph-Env [12]	2	-	×	✓	×	×
RLOR [116]	2	2	×	✓	✓	×
RoutingArena [111]	1	8	✓	×	×	×
Jumanji [14]	22	3	✓	✓	×	×
RL4CO (ours)	27 [‡]	23	✓	✓	✓	✓

[†] We consider as *baselines* ad-hoc network architectures (i.e., policies) and RL algorithms from the literature.

[‡] We also consider the possible 16 combinations of environments generated by the unified Multi-Task VRP, as they have been historically considered separate environments in the NCO literature.

75 **Related Benchmark Libraries.** Despite the variety of general-purpose RL software libraries [18,
76 70, 96, 119, 24, 33, 81], there is a lack of a unified and extensive benchmark for CO problems. Balaji
77 et al. [4] propose an RL benchmark for Operations Research (OR) with a PPO baseline [100]; Hubbs
78 et al. [42], Biagioni et al. [12] provide a collection of OR environments. Wan et al. [116] propose
79 a general-purpose library for OR, and benchmarks the canonical TSP and CVRP environments.
80 However, a major downside of the above libraries is that they cannot be massively parallelized due
81 to their reliance on the OpenAI Gym API, which can only run on CPU, unlike RL4CO, which is
82 based on the TorchRL [15], a recent official PyTorch [92] library for RL that enables hardware-
83 accelerated execution of both environments and algorithms. Prouvost et al. [94] introduces a library
84 specialized for CO problems that work in combination with traditional MILP [71] solvers. We also
85 mention Routing Arena [111], whose scope is different from RL4CO, namely, comparing NCO and
86 classical solvers only for the CVRP. The most related work is Jumanji [14], which provides a variety
87 of CO environments written in JAX [16] that can be hardware-accelerated alongside an actor-critic
88 baseline. While Jumanji is an RL environment suite, RL4CO is a full-stack library that integrates
89 environments, policies, RL algorithms under a unified framework.

90 3 RL4CO: Taxonomy

91 We describe the RL4CO taxonomy, categorizing components into *Environments*, *Policies*, and *RL*
92 *Algorithms*. Then we translate the taxonomy to implementation in § 4.

93 **Environments.** Given a CO problem instance \mathbf{x} , we formulate the solution-generating procedure as
94 a Markov Decision Process (MDP) characterized by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$ as follows. **State** \mathcal{S}
95 is the space of states that represent the given problem \mathbf{x} and the current partial solution being updated
96 in the MDP. **Action** \mathcal{A} is the action space, which includes all feasible actions a_t that can be taken at
97 each step t . **State Transition** \mathcal{T} is the deterministic state transition function $s_{t+1} = \mathcal{T}(s_t, a_t)$ that
98 updates a state s_t to the next state s_{t+1} . **Reward** \mathcal{R} is the reward function $\mathcal{R}(s_t, a_t)$ representing
99 the immediate reward received after taking action a_t in state s_t . Finally, $\gamma \in [0, 1]$ is a discount
100 factor that determines the importance of future rewards. Since the state transition is deterministic,
101 we represent the solution for a problem \mathbf{x} as a sequence of T actions $\mathbf{a} = (a_1, \dots, a_T)$. Then the
102 total return $\sum_{t=1}^T \mathcal{R}(s_t, a_t)$ translates to the negative cost function of the CO problem.

103 **Policies.** The policies can be categorized into constructive policies, which generate a solution from
104 scratch, and improvement policies, which refine an existing solution.

105 *Constructive policies.* A policy π is used to construct a solution from scratch for a given problem
106 instance \mathbf{x} . It can be further categorized into autoregressive (AR) and non-autoregressive (NAR)
107 policies. An AR policy is composed by an encoder f that maps the instance \mathbf{x} into an embedding
108 space $\mathbf{h} = f(\mathbf{x})$ and by a decoder g that iteratively determines a sequence of actions \mathbf{a} as follows:

$$a_t \sim g(a_t | a_{t-1}, \dots, a_0, s_t, \mathbf{h}), \quad \pi(\mathbf{a} | \mathbf{x}) \triangleq \prod_{t=1}^{T-1} g(a_t | a_{t-1}, \dots, a_0, s_t, \mathbf{h}). \quad (1)$$

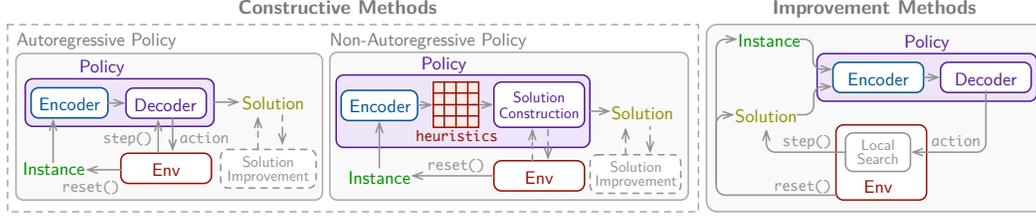


Figure 1: Overview of different types of policies and their modularization in RL4CO.

109 A NAR policy encodes a problem \mathbf{x} into a heuristic $\mathcal{H} = f(\mathbf{x}) \in \mathbb{R}_+^N$, where N is the number of
 110 possible assignments across all decision variables. Each number in \mathcal{H} represents a (unnormalized)
 111 probability of a particular assignment. To obtain a solution \mathbf{a} from \mathcal{H} , one can sample a sequence
 112 of assignments from \mathcal{H} while dynamically masking infeasible assignments to meet problem-specific
 113 constraints. It can also guide a search process, e.g., Ant Colony Optimization [28, 125, 56], or be
 114 incorporated into hybrid frameworks [127]. Here, the heuristic helps identify promising transitions
 115 and improve the efficiency of finding an optimal or near-optimal solution.

116 *Improvement policies.* A policy can be used for improving an initial solution $\mathbf{a}^0 = (a_0^0, \dots, a_{T-1}^0)$
 117 into another one potentially with higher quality, which can be formulated as follows:

$$\mathbf{a}^k \sim g(\mathbf{a}^0, \mathbf{h}), \quad \pi(\mathbf{a}^k | \mathbf{a}^0, \mathbf{x}) \triangleq \prod_{k=1}^{K-1} g(\mathbf{a}^k | \mathbf{a}^{k-1}, \dots, \mathbf{a}^0, \mathbf{h}), \quad (2)$$

118 where \mathbf{a}^k is the k -th updated solution and K is the budget for number of improvements. This process
 119 allows continuous refinement for a long time to enhance the solution quality.

120 **RL Algorithms.** The RL objective is to learn a policy π that maximizes the expected cumulative
 121 reward (or equivalently minimizes the cost) over the distribution of problem instances:

$$\theta^* = \operatorname{argmax}_{\theta} \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} \left[\mathbb{E}_{\pi(\mathbf{a}|\mathbf{x})} \left[\sum_{t=0}^{T-1} \gamma^t \mathcal{R}(s_t, a_t) \right] \right], \quad (3)$$

122 where θ is the set of parameters of π and $P(\mathbf{x})$ is the distribution of problem instances. Eq. (3) can
 123 be solved using algorithms such as variations of REINFORCE [109], Advantage Actor-Critic (A2C)
 124 methods [59], or Proximal Policy Optimization (PPO) [100]. These algorithms are employed to train
 125 the policy network π , by transforming the maximization problem in Eq. (3) into a minimization
 126 problem involving a loss function, which is then optimized using gradient descent algorithms. For
 127 instance, the REINFORCE loss function gradient is given by:

$$\nabla_{\theta} \mathcal{L}_a(\theta | \mathbf{x}) = \mathbb{E}_{\pi(\mathbf{a}|\mathbf{x})} [(R(\mathbf{a}, \mathbf{x}) - b(\mathbf{x})) \nabla_{\theta} \log \pi(\mathbf{a} | \mathbf{x})], \quad (4)$$

128 where $b(\cdot)$ is a baseline function used to stabilize training and reduce gradient variance. We also
 129 distinguish between two types of RL (pre)training: 1) *inductive* and 2) *transductive* RL. In inductive
 130 RL, the focus is on learning patterns from the training dataset to generalize to new instances, thus
 131 amortizing the inference procedure. Conversely, transductive RL (or test-time optimization) opti-
 132 mizes parameters during testing on target instances. Typically, a policy π is trained using inductive
 133 RL, followed by transductive RL for test-time optimization.

134 4 RL4CO: Library Structure

135 RL4CO is a unified reinforcement learning (RL) for Combinatorial Optimization (CO) library that
 136 aims to provide a *modular, flexible, and unified* code base for training and evaluating RL for CO
 137 methods with extensive benchmarking capabilities on various settings. As shown in Fig. 2, RL4CO
 138 decouples the major components of an RL pipeline, prioritizing their reusability in the implementa-
 139 tion. Following also the taxonomy of § 3, the main components are: (§ 4.1) Environments, (§ 4.2)
 140 Policies, (§ 4.3) RL algorithms, (§ 4.4) Utilities, and (§ 4.5) Environments & Baselines Zoo.

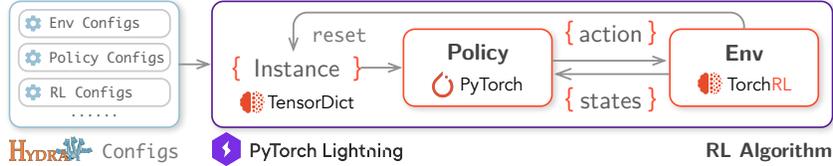


Figure 2: Overview of the RL4CO pipeline: from configurations to training a policy.

141 4.1 Environments

142 Environments in RL4CO fully specify the CO problems and their logic. They are based on the
 143 RL4COEnvBase class that extends from the EnvBase in TorchRL [15]. A modular generator
 144 can be provided to the environment. The generator provides CO instances to the environ-
 145 ment, and different generators can be used to generate different data distributions. Static in-
 146 stance data and dynamic variables, such as the current state s_t , current solution a^k for im-
 147 provement environments, policy actions a_t , rewards, and additional information are passed in a
 148 *stateless* fashion in a TensorDict [86], that we call `td`, through the environment `reset` and
 149 `step` functions. Additionally, our environment API contains several functions, such as `render`,
 150 `check_solution_validity`, `select_start_nodes` (i.e., for POMO-based optimization [62])
 151 and optional API as `local_search` solution improvement.

152 It is noteworthy that RL4CO enhances the efficiency of environments when compared to vanilla
 153 TorchRL, by overriding and optimizing some methods in TorchRL EnvBase. For instance, our new
 154 `step` method brings a decrease of up to 50% in latency and halves the memory impact by avoiding
 155 saving duplicate components in the stateless TensorDict.

156 4.2 Policies

157 Policies in RL4CO are subclasses of PyTorch’s `nn.Module` and contain the encoding-decoding
 158 logic and neural network parameters θ . Different policies in the RL4CO "zoo" can inherit from
 159 metaclasses like `ConstructivePolicy` or `ImprovementPolicy`. We modularize components to
 160 process raw features into the embedding space via a parametrized function ϕ_ω , called *feature em-*
 161 *beddings*. 1) *Node Embeddings* ϕ_n : transform m_n node features of instances x from the feature
 162 space to the embedding space h , i.e., $[B, N, m_n] \rightarrow [B, N, h]$. 2) *Edge Embeddings* ϕ_e : trans-
 163 form m_e edge features of instances x from the feature space to the embedding space h , i.e.,
 164 $[B, E, m_e] \rightarrow [B, E, h]$, where E is the number of edges. 3) *Context Embeddings* ϕ_c : capture
 165 contextual information by transforming m_c context features from the current decoding step s_t from
 166 the feature space to the embedding space h , i.e., $[B, m_c] \rightarrow [B, h]$, for nodes or edges. Overall,
 167 Fig. 3 illustrates a generic constructive AR policy in RL4CO, where the feature embeddings are ap-
 168 plied similarly to other types of policies. Embeddings can be automatically selected by RL4CO at
 169 runtime by simply passing the `env_name` to the policy. Additionally, we allow for granular control
 170 of any higher-level policy component independently, such as encoders and decoders.

171 4.3 RL Algorithms

172 RL algorithms in RL4CO define the process that takes the Environment with its problem in-
 173 stances and the Policy to optimize its parameters θ . The parent class of algorithms is the
 174 RL4COLitModule, inheriting from PyTorch Lightning’s `pl.LightningModule` [31]. This al-
 175 lows for granular support of various methods including the `[train, val, test]_step`, auto-
 176 matic logging with several logging services such as Wandb via `log_metrics`, automatic optimizer
 177 configuration via `configure_optimizers` and several useful callbacks for RL methods such as
 178 `on_train_epoch_end`. RL algorithms are additionally attached to an RL4COTrainer, a wrap-
 179 per we made with additional optimizations around `pl.Trainer`. This module seamlessly supports
 180 features of modern training pipelines, including logging, checkpoint management, mixed-precision
 181 training, various hardware acceleration supports (e.g., CPU, GPU, TPU, and Apple Silicon), and
 182 multi-device hardware accelerator in distributed settings [69]. For instance, using mixed-precision

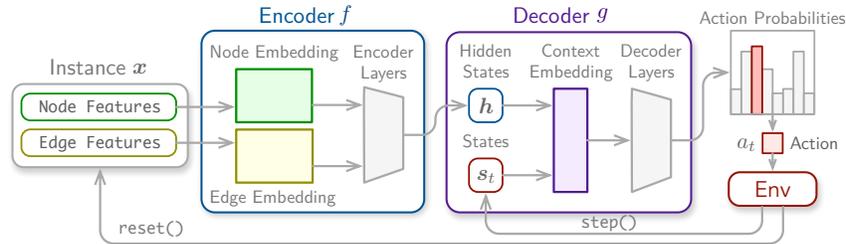


Figure 3: Overview of modularized RL4CO policies. Any component such as the encoder/decoder structure and feature embeddings can be replaced and thus the model is adaptable to various new environments.

183 training significantly decreases training time without sacrificing much convergence and enables us
 184 to leverage recent routines, e.g., FlashAttention [26, 25], which we investigate in Appendix.

185 4.4 Utilities

186 **Configuration Management.** Optionally, but usefully, we adopt Hydra [123], an open-source
 187 Python framework that enables hierarchical config management, making it easier to manage com-
 188 plex configurations and experiments with different settings as shown in Appendix. Hydra addi-
 189 tionally allows for automatically parsing parameters (un-)defined in configs - i.e., `python run.py`
 190 `experiment=routing/pomo env=cvrp env.generator_params.num_loc=50` launches an
 191 experiment defined under `routing/pomo` and changes the environment to CVRP with 50 locations.

192 **Decoding Schemes.** Decoding schemes handle the logic of model logits z by applying preprocess-
 193 ing, such as masking of infeasible actions and/or additional techniques to select better actions during
 194 training and testing. We implement the model and problem-agnostic decoding schemes under the
 195 `DecodingStrategy` class in the RL4CO codebase that can be easily reused: 1) *Greedy*, which
 196 selects the action with the highest probability; 2) *Sampling*, which samples `n_samples` solutions
 197 from the current masked probability distribution of the policy, incorporating sampling strategies like
 198 2.a) Softmax Temperature τ , 2.b) top-k sampling [61], and 2.c) top-p (or Nucleus) sampling [38]
 199 (more details in Appendix); 3) *Multistart*, which enforces diverse starting actions as demonstrated in
 200 POMO [62], such as starting from different cities in the Traveling Salesman Problem (TSP) with N
 201 nodes; 4) *Augmentation*, which applies transformations to instances, such as random rotations and
 202 flipping in Euclidean problems [55], to create an augmented set of problems.

203 **Documentation, Tutorials, and Testing.** We release extensive documentation to make it as acces-
 204 sible as possible for both newcomers and experts. RL4CO can be easily installed by running `pip`
 205 `install rl4co` with open-source code available at <https://github.com/ai4co/rl4co>. Sev-
 206 eral tutorials and examples are also available under the `examples/` folder. We thoroughly test our
 207 library via continuous integration on multiple Python versions and operating systems. The following
 208 code snippet shows minimalistic code that can train a model in a few lines:

```

from rl4co.envs.routing import TSPEnv, TSPGenerator
from rl4co.models import AttentionModelPolicy, POMO
from rl4co.utils import RL4COTrainer
# Instantiate generator and environment
generator = TSPGenerator(num_loc=50, loc_distribution="uniform")
env = TSPEnv(generator)
# Create policy and RL model
policy = AttentionModelPolicy(env_name=env.name, num_encoder_layers=6)
model = POMO(env, policy, batch_size=64)
# Instantiate Trainer and fit
trainer = RL4COTrainer(max_epochs=10, accelerator="gpu", precision="16-mixed")
trainer.fit(model)

```

209 4.5 Environments & Baselines Zoo

210 **Environments.** We include benchmarking from the following environments, divided into four areas. 1) **Routing:** Traveling Salesman Problem (TSP) [65], Capacitated Vehicle Routing Problem (CVRP) [13], Orienteering Problem (OP) [64, 21], Prize Collecting TSP (PCTSP) [5], Pickup and Delivery Problem (PDP) [50, 99] and Multi-Task VRP (MTVRP) [72, 131, 9] (which modularizes with 16 problem variants including the basic VRPTW, OVRP, VRPB, VRPL and VRPs with their constraint combinations); 2) **Scheduling:** Flexible Job Shop Scheduling Problem (FJSSP) [17], Job Shop Scheduling Problem (JSSP) [97] and Flow Shop Scheduling Problem (FJSP); 3) **Electronic Design Automation:** multiple Decap Placement Problem (mDPP) [53]; 4) **Graph:** Facility Location Problem (FLP) [30] and Max Cover Problem (MCP) [51].

219 **Baseline Zoo.** Given that several works contribute to both new policies and new RL algorithm variations, we list the papers we reproduce. For 1) **Constructive AR** methods, we include the Attention Model (AM) [60], Ptr-Net [115], POMO [62], MatNet [63], HAM [67], SymNCO [55], PolyNet [41], MTPOMO [72], MVMoE [131], L2D [128], HGNN [106] and DevFormer [53]. For 2) **Constructive NAR** methods, we benchmark Ant Colony Optimization-based DeepACO [125] and GFACS [56] as well as the hybrid NAR/AR GLOP [127]. 3) **Improvement methods** include DACT [78], N2S [79] and NeuOpt [80]. We also include 4) **General-purpose RL** algorithm from the literature, including REINFORCE [109] with various baselines, Advantage Actor-Critic (A2C) [59] and Proximal Policy Optimization (PPO) [100] that can be readily be combined with any policy. 228 Finally, we include 5) **Active search** (i.e., Transductive RL) methods AS [6] and EAS [40].

229 5 Benchmarking Study

230 We perform several benchmarking studies with our unified RL4CO library. Given the limited space, 231 we invite the reader to check out the [Appendix](#) for supplementary material.

232 5.1 Flexibility and Modularity

233 **Changing policy components.** The integration of many state-of-the-art methods in RL4CO from 234 the NCO field in a modular framework makes it easy to implement and improve upon state-of-the- 235 art neural solvers for complex CO problems with only a few lines of code and improve upon them.²

236 We demonstrate this in [Table 2](#) for the FJSSP by 237 gradually replacing or adding elements to the 238 original SotA policy [106]. First, replacing the 239 HGNN encoder with the more expressive Mat- 240 Net encoder [63] already improves the average 241 makespan by around 7%. Further improvements 242 can be achieved by replacing the MLP 243 decoder with the Pointer mechanism in the AM 244 decoder [60] with gaps to BKS around $3\times$ 245 lower compared to the original policy in Song 246 et al. [106] even with greedy performance.

Table 2: Solutions obtained with RL4CO for the FJSSP with different model configurations.

Encoder / Decoder		FJSSP	
		10×5	20×5
HGNN + MLP (g.) [106]	Obj.	111.82	211.21
	Gap	15.8%	12.1%
MatNet + MLP (g.)	Obj.	103.91	197.92
	Gap	7.6%	5.0%
MatNet + Pointer (g.)	Obj.	101.17	196.3
	Gap	4.8%	4.2%
MatNet + Pointer (s. x128)	Obj.	98.31	192.02
	Gap	1.8%	1.9%

247 5.2 Constructive Policies

248 **Mind Your Baseline.** In on-policy RL, which is often employed in RL4CO due to fast reward 249 function evaluations, several different REINFORCE baselines have been proposed to improve the 250 performance. We benchmark several RL algorithms training constructive policies for routing prob- 251 lems of node size 50, whose underlying architecture is based on the encoder-decoder Attention 252 Model [60] and whose main difference lies in how the REINFORCE baseline is calculated (we ad- 253 ditionally train the AM with PPO as further reference). For a fair comparison, we run all baselines

²The different model configurations shown here can be obtained by simply changing the Hydra configura- tion file like the one shown in [Appendix](#).

254 in controlled settings with the same number of optimization steps and report results in Table 3. We
 255 note that A2C generally underperforms other baselines. Such performance can be attributed to the
 256 fact that since in routing problems, the rewards are sparse (i.e., can only be calculated upon solving
 257 an entire problem), estimating the value of an entire instance x is inherently a challenging task.
 258 Interestingly, while POMO [62], which takes as a baseline the shared baseline of all routes
 259 forcing each starting node to be different, may work well as baselines for problems in which
 260 near-optimal solutions can be constructed from any node (e.g., TSP), this may not be true for
 261 other problems such as the Orienteering Problem (OP): the reason is that in OP only a *subset*
 262 of nodes should be selected in an optimal solution, while several states will be discarded. Hence,
 263 forcing the policy to select all of them makes up for a poor baseline. We remark that while SymNCO
 264 (whose shared baseline involves symmetric rotations and flips) [55] may perform well in Euclidean
 265 problems, this is not applicable in non-Euclidean CO, including asymmetric routing problems and
 266 scheduling. We found similar trends regarding actor-critic methods as A2C and PPO in the EDA
 267 mDPP problem [53], which we report in Appendix. Namely, a greedy rollout baseline [60] can do
 268 better than value-based methods due to the challenging task of instance value estimation.

Table 3: Optimality Gap obtained via greedy decoding.

Method	TSP	CVRP	OP	PCTSP	PDP
A2C	2.22	7.09	8.64	14.96	10.02
AM-Rollout	1.41	5.30	4.40	2.46	9.88
POMO	0.89	3.99	14.26	11.61	10.64
Sym-NCO	0.47	4.61	3.09	2.12	7.73
AM-PPO	0.92	4.60	3.05	2.45	8.31

273 **Decoding Schemes.** The solution quality of NCO
 274 solvers often shows significant improvements in per-
 275 formance to different decoding schemes, even with
 276 the exact NCO solvers. We evaluate the trained
 277 solver with different decoding schemes and settings
 278 as shown in Fig. 4.

279 **Generalization.** Using RL4CO, we can easily eval-
 280 uate the generalization performance of existing base-
 281 lines by employing supported environments that in-
 282 corporate various VRP variant tasks and instance
 283 distributions (termed MTPOMO and MDPOMO, re-
 284 spectively). Empirical results on CVRPLib, shown
 285 in Table 4, reveal that training on different tasks
 286 significantly enhances generalization performance.
 287 This finding underscores the necessity of building foundational models across diverse CO domains.

288 **Large-Scale Instances.** We evaluate large-scale CVRP instances of thousands of nodes, with more
 289 visualizations and scaling in Appendix. The last row of Table 5 illustrates the performance of the
 290 hybrid NAR/AR GLOP [127], while others refer to reproduced results from Ye et al. [127]. Our
 291 implementation in RL4CO improves the performance in not only speed but also solution quality.

292 5.3 Combining Construction and Improvement: Best of Both Worlds?

293 While constructive policies can build solutions in seconds, their performance is often limited, even
 294 with advanced decoding schemes such as sampling or augmentations. On the other hand, improve-
 295 ment methods are more suitable for larger computing budgets. We benchmark models on TSP with
 296 50 nodes: the AR constructive method POMO [62] and the improvement methods DACT [78] and
 297 NeuOpt [80]. In the original implementation, DACT and NeuOpt started from a solution constructed
 298 randomly. To further demonstrate the flexibility of RL4CO, we show that bootstrapping improve-
 299 ment methods with constructive ones enhance convergence speed. Fig. 5 shows that bootstrapping
 300 with a pre-trained POMO policy significantly enhances the convergence speed. To further investi-
 301 gate the performance, we report the Primal Integral (PI) [8, 113, 111], which evaluates the evolution
 302 of solution quality over time. Improvement methods alone, such as DACT and NeuOpt, achieve 2.99
 303 and 2.26 respectively, while sampling from POMO achieves 0.08. This shows that the “area under
 304 the curve” can be better even if the final solution is worse for constructive methods. Bootstrapping

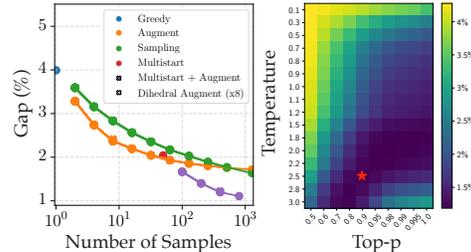


Figure 4: Decoding schemes study of POMO on CVRP50. [Left]: Pareto front of decoding schemes by the number of samples; [Right]: performance of sampling with different temperatures τ and p values for top- p sampling.

Benchmark	POMO		MTPOMO		MDPOMO	
	Obj.	Gap	Obj.	Gap	Obj.	Gap
Set A	1075	3.13%	1076	3.20%	1074	2.97%
Set B	996	3.41%	1003	4.06%	995	3.26%
Set E	761	5.04%	760	4.82%	762	5.07%
Set F	813	13.52%	798	12.09%	825	13.66%
Set M	1259	16.37%	1234	13.58%	1263	16.03%
Set P	620	6.72%	608	3.72%	613	5.04%
Set X	73953	16.80%	73763	16.69%	81848	23.69%

Table 4: Results on CVRPLIB instances with models trained on $N = 50$. Greedy multi-start decoding is used.

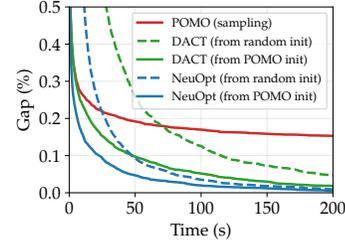


Figure 5: Bootstrapping improvement with constructive methods.

305 with POMO then improves DACT and NeuOpt to 0.08 and 0.04 respectively, showing the benefits
 306 of modularity and hybridization of different components.

307 6 Discussion

308 6.1 Limitations and Future Directions

309 While RL4CO is an efficient and modular library specialized in CO problems, it might not
 310 be suitable for any other task due to a number of area-specific optimizations, and we do not
 311 expect it to seamlessly integrate with, for instance, OpenAI Gym wrappers without some
 312 modifications. Another limitation of the library is its scope so far, namely RL. In fact, extending
 313 the library to support supervised methods and creating a comprehensive "AI4CO" library
 314 could benefit the whole NCO community. We additionally identify in Foundation Models³ for CO and related scalable architectures a promising
 315 area of future research to overcome generalization issues across tasks and distributions, for which
 316 we provided some early clues.

Table 5: Performance on large-scale CVRP instances.

	CVRP1K		CVRP2K		CVRP7K	
	Obj.	Time (s)	Obj.	Time (s)	Obj.	Time (s)
LKH-3	46.4	6.2	64.9	20	245.0	501
AM	61.4	0.6	114.4	1.9	354.3	26
TAM(AM)	50.1	0.8	74.3	2.2	233.4	26
TAM(LKH-3)	46.3	1.8	64.8	5.6	196.9	33
GLOP-G(AM)*	47.1	0.4	63.5	1.2	191.7	2.4
GLOP-G(LKH-3)*	45.9	1.1	63.0	1.5	191.2	5.8
GLOP-G(AM)	46.9	0.3	64.7	0.7	190.9	2.0
GLOP-G(LKH-3)	45.5	0.5	62.8	0.8	190.1	3.9

323 6.2 Long-term Plans

324 Our long-term plan is to become the go-to RL for CO benchmark library. While not strictly tied
 325 to implementation and benchmarking, we are committed to helping resolve issues and questions
 326 from the community. For this purpose, we created a Slack workspace (link available in the online
 327 documentation) that by now has attracted more than 130 researchers. It is our hope that our work
 328 will ultimately benefit the NCO field with new ideas and collaborations.

329 7 Conclusion

330 This paper introduces RL4CO, a modular, flexible, and unified Reinforcement Learning (RL) for
 331 Combinatorial Optimization (CO) benchmark. We provide a comprehensive taxonomy from envi-
 332 ronments to policies and RL algorithms that translate from theory to practice to software level. Our
 333 benchmark library aims to fill the gap in unifying implementations in RL for CO by utilizing sev-
 334 eral best practices with the goal of providing researchers and practitioners with a flexible starting
 335 point for NCO research. We provide several experimental results with insights and discussions that
 336 can help identify promising research directions. We hope that our open-source library will provide
 337 a solid starting point for NCO researchers to explore new avenues and drive advancements. We
 338 warmly welcome researchers and practitioners to actively participate and contribute to RL4CO.

³<https://github.com/ai4co/awesome-fm4co>

339 Acknowledgements

340 We want to express our gratitude towards anonymous reviewers of previous submissions who greatly
341 helped us improve our paper. Even though rejections were not easy at first, they helped us refine
342 our benchmark. Importantly, through our journey, we got to know several outstanding researchers
343 in the community, who gave us even more motivation and meaning behind our work. We would
344 also like to thank people in the AI4CO open research community who have contributed, and those
345 who will, to RL4CO. We also thank OMELET for supporting us with additional compute. We
346 invite practitioners and researchers to join us and contribute with bug reporting, feature requests, or
347 collaboration ideas. A special thanks also goes to the TorchRL team for helping us in solving issues
348 and improving the library.

349 Potential Broader Impact

350 This paper presents work in the field of AI4CO. The main consequence may be that AI methods to
351 solve CO problems may become accessible to the broad public, as our library is open source and
352 readily available on GitHub. We do not see potential negative societal consequences as of today.

353 Funding

354 This work was supported by a grant of the KAIST-KT joint research project through AI2XL Labo-
355 ratory, Institute of Convergence Technology, funded by KT [Project No. G01210696, Development
356 of Multi-Agent Reinforcement Learning Algorithm for Efficient Operation of Complex Distributed
357 Systems] and by the Institute of Information & communications Technology Planning & Evaluation
358 (IITP) grant funded by the Korean government(MSIT)[2022-0-01032, Development of Collective
359 Collaboration Intelligence Framework for Internet of Autonomous Things].

360 References

- 361 [1] A. AhmadiTeshnizi, W. Gao, and M. Udell. OptiMUS: Scalable optimization modeling with
362 (mi) lp solvers and large language models. In *International Conference on Machine Learning*,
363 2024.
- 364 [2] S. Ahn, Y. Seo, and J. Shin. Learning what to defer for maximum independent sets. In
365 *International Conference on Machine Learning*, pages 134–144. PMLR, 2020.
- 366 [3] K. Ali, W. Alsalih, and H. Hassanein. Set-cover approximation algorithms for load-aware
367 readers placement in RFID networks. In *2011 IEEE international conference on communica-*
368 *tions (ICC)*, pages 1–6. IEEE, 2011.
- 369 [4] B. Balaji, J. Bell-Masterson, E. Bilgin, A. Damianou, P. M. Garcia, A. Jain, R. Luo, A. Mag-
370 giar, B. Narayanaswamy, and C. Ye. Orl: Reinforcement learning benchmarks for online
371 stochastic optimization problems. *arXiv preprint arXiv:1911.10641*, 2019.
- 372 [5] E. Balas. The prize collecting traveling salesman problem. *Networks*, 19(6):621–636, 1989.
- 373 [6] I. Bello, H. Pham, Q. V. Le, M. Norouzi, and S. Bengio. Neural combinatorial optimization
374 with reinforcement learning, 2017.
- 375 [7] Y. Bengio, A. Lodi, and A. Prouvost. Machine learning for combinatorial optimization: a
376 methodological tour d’horizon. *European Journal of Operational Research*, 290(2):405–421,
377 2021.
- 378 [8] T. Berthold. Measuring the impact of primal heuristics. *Operations Research Letters*, 41(6):
379 611–614, 2013.

- 380 [9] F. Berto, C. Hua, N. G. Zepeda, A. Hottung, N. Wouda, L. Lan, K. Tierney, and J. Park.
381 RouteFinder: Towards foundation models for vehicle routing problems, 2024. GitHub repos-
382 itory: <https://github.com/ai4co/routefinder>.
- 383 [10] K. Bestuzheva, M. Besançon, W.-K. Chen, A. Chmiela, T. Donkiewicz, J. van Doornmalen,
384 L. Eifler, O. Gaul, G. Gamrath, A. Gleixner, et al. The SCIP optimization suite 8.0. *arXiv*
385 *2112.08872*, 2021.
- 386 [11] J. Bi, Y. Ma, J. Wang, Z. Cao, J. Chen, Y. Sun, and Y. M. Chee. Learning generalizable models
387 for vehicle routing problems via knowledge distillation. *Advances in Neural Information*
388 *Processing Systems*, 35:31226–31238, 2022.
- 389 [12] D. Biagioni, C. E. Tripp, S. Clark, D. Duplyakin, J. Law, and P. C. S. John. graphenv: a
390 python library for reinforcement learning on graph search spaces. *Journal of Open Source*
391 *Software*, 7(77):4621, 2022.
- 392 [13] L. Bodin. Routing and scheduling of vehicles and crews. *Computer & Operations Research*,
393 10(2):69–211, 1983.
- 394 [14] C. Bonnet, D. Luo, D. Byrne, S. Surana, S. Abramowitz, P. Duckworth, V. Coyette,
395 L. I. Midgley, E. Tegegn, T. Kalloniatis, O. Mahjoub, M. Macfarlane, A. P. Smit,
396 N. Grinsztajn, R. Boige, C. N. Waters, M. A. Mimouni, U. A. M. Sob, R. de Kock,
397 S. Singh, D. Furelos-Blanco, V. Le, A. Pretorius, and A. Laterre. Jumanji: a
398 diverse suite of scalable reinforcement learning environments in jax, 2024. URL
399 [InternationalConferenceonLearningRepresentations](https://arxiv.org/abs/2406.18250).
- 400 [15] A. Bou, M. Bettini, S. Dittert, V. Kumar, S. Sodhani, X. Yang, G. D. Fabritiis, and V. Moens.
401 TorchRL: A data-driven decision-making library for pytorch. In *International conference on*
402 *learning representations*, 2024.
- 403 [16] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula,
404 A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable trans-
405 formations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- 406 [17] P. Brandimarte. Routing and scheduling in a flexible job shop by tabu search. *Annals of*
407 *Operations research*, 41(3):157–183, 1993.
- 408 [18] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.
409 Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 410 [19] S. Brody, U. Alon, and E. Yahav. How attentive are graph attention networks? In *Interna-*
411 *tional Conference on Learning Representations*, 2019.
- 412 [20] F. Bu, H. Jo, S. Y. Lee, S. Ahn, and K. Shin. Tackling prevalent conditions in unsupervised
413 combinatorial optimization: Cardinality, minimum, covering, and more. In *International*
414 *Conference on Machine Learning*, 2024.
- 415 [21] I.-M. Chao, B. L. Golden, and E. A. Wasil. A fast and effective heuristic for the orienteering
416 problem. *European journal of operational research*, 88(3):475–489, 1996.
- 417 [22] X. Chen and Y. Tian. Learning to perform local rewriting for combinatorial optimization. In
418 *Advances in Neural Information Processing Systems*, 2019.
- 419 [23] H. Dai, E. B. Khalil, Y. Zhang, B. Dilkina, and L. Song. Learning combinatorial optimization
420 algorithms over graphs. In *Advances in Neural Information Processing Systems*, volume 30,
421 2017.
- 422 [24] S. Dalton et al. Accelerating reinforcement learning through gpu atari emulation. *Advances*
423 *in Neural Information Processing Systems*, 33:19773–19782, 2020.

- 424 [25] T. Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning.
425 *arXiv preprint arXiv:2307.08691*, 2023.
- 426 [26] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré. Flashattention: Fast and memory-efficient
427 exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:
428 16344–16359, 2022.
- 429 [27] V. C. David Applegate, Robert Bixby and W. Cook. Concorde TSP solver, 2023. URL
430 <https://www.math.uwaterloo.ca/tsp/concorde/index.html>.
- 431 [28] M. Dorigo and T. Stützle. *Ant colony optimization: overview and recent advances*. Springer,
432 2019.
- 433 [29] D. Drakulic, S. Michel, F. Mai, A. Sors, and J.-M. Andreoli. Bq-nco: Bisimulation quoti-
434 enting for generalizable neural combinatorial optimization. *Advances in Neural Information*
435 *Processing Systems*, 2023.
- 436 [30] Z. Drezner and H. W. Hamacher. *Facility location: applications and theory*. Springer Science
437 & Business Media, 2004.
- 438 [31] W. Falcon and The PyTorch Lightning team. PyTorch Lightning, 3 2019. URL [https://](https://github.com/Lightning-AI/lightning)
439 github.com/Lightning-AI/lightning.
- 440 [32] M. Fischetti, J. J. S. Gonzalez, and P. Toth. Solving the orienteering problem through branch-
441 and-cut. *INFORMS Journal on Computing*, 10(2):133–148, 1998.
- 442 [33] C. D. Freeman, E. Frey, A. Raichuk, S. Girgin, I. Mordatch, and O. Bachem. Brax - a
443 differentiable physics engine for large scale rigid body simulation, 2021. URL [http://](http://github.com/google/brax)
444 github.com/google/brax.
- 445 [34] B. L. Golden, L. Levy, and R. Vohra. The orienteering problem. *Naval Research Logistics*
446 (*NRL*), 34(3):307–318, 1987.
- 447 [35] L. Gurobi Optimization. Gurobi optimizer reference manual, 2021. URL [http://www.](http://www.gurobi.com)
448 [gurobi.com](http://www.gurobi.com).
- 449 [36] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs.
450 *Advances in neural information processing systems*, 30, 2017.
- 451 [37] K. Helsgaun. An extension of the lin-kernighan-helsgaun tsp solver for constrained traveling
452 salesman and vehicle routing problems. *Roskilde: Roskilde University*, 12 2017. doi: 10.
453 13140/RG.2.2.25569.40807.
- 454 [38] A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi. The curious case of neural text degener-
455 ation. *arXiv preprint arXiv:1904.09751*, 2019.
- 456 [39] A. Hottung, B. Bhandari, and K. Tierney. Learning a latent search space for routing problems
457 using variational autoencoders. In *International Conference on Learning Representations*,
458 2020.
- 459 [40] A. Hottung, Y.-D. Kwon, and K. Tierney. Efficient active search for combinatorial optimiza-
460 tion problems. *International conference on learning representations*, 2022.
- 461 [41] A. Hottung, M. Mahajan, and K. Tierney. PolyNet: Learning diverse solution strategies for
462 neural combinatorial optimization. *arXiv preprint arXiv:2402.14048*, 2024.
- 463 [42] C. D. Hubbs, H. D. Perez, O. Sarwar, N. V. Sahinidis, I. E. Grossmann, and J. M. Wassick.
464 Or-gym: A reinforcement learning library for operations research problems. *arXiv preprint*
465 *arXiv:2008.06319*, 2020.

- 466 [43] J. Hwang, J. S. Pak, D. Yoon, H. Lee, J. Jeong, Y. Heo, and I. Kim. Enhancing on-die pdn
467 for optimal use of package pdn with decoupling capacitor. In *2021 IEEE 71st Electronic
468 Components and Technology Conference (ECTC)*, pages 1825–1830, 2021. doi: 10.1109/
469 ECTC32696.2021.00288.
- 470 [44] Z. Iklassov, Y. Du, F. Akimov, and M. Takac. Self-guiding exploration for combinatorial
471 problems. *arXiv preprint arXiv:2405.17950*, 2024.
- 472 [45] L. Ivan. Capacitated vehicle routing problem library. [http://vrp.atd-lab.inf.puc-
473 rio.br/index.php/en/](http://vrp.atd-lab.inf.puc-rio.br/index.php/en/). 2014.
- 474 [46] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton. Adaptive mixtures of local experts.
475 *Neural computation*, 3(1):79–87, 1991.
- 476 [47] Y. Jiang, Y. Wu, Z. Cao, and J. Zhang. Learning to solve routing problems via distributionally
477 robust optimization. In *36th AAAI Conference on Artificial Intelligence*, 2022.
- 478 [48] C. K. Joshi, T. Laurent, and X. Bresson. An efficient graph convolutional network technique
479 for the travelling salesman problem. *arXiv preprint arXiv:1906.01227*, 2019.
- 480 [49] J. Juang, L. Zhang, Z. Kiguradze, B. Pu, S. Jin, and C. Hwang. A modified genetic algorithm
481 for the selection of decoupling capacitors in pdn design. In *2021 IEEE International Joint
482 EMC/SI/PI and EMC Europe Symposium*, pages 712–717, 2021. doi: 10.1109/EMC/SI/PI/
483 EMCEurope52599.2021.9559292.
- 484 [50] B. Kalantari, A. V. Hill, and S. R. Arora. An algorithm for the traveling salesman problem
485 with pickup and delivery customers. *European Journal of Operational Research*, 22(3):377–
486 386, 1985.
- 487 [51] S. Khuller, A. Moss, and J. S. Naor. The budgeted maximum coverage problem. *Information
488 processing letters*, 70(1):39–45, 1999.
- 489 [52] D. Kikuta, H. Ikeuchi, K. Tajiri, and Y. Nakano. RouteExplainer: An explanation framework
490 for vehicle routing problem. In *Pacific-Asia Conference on Knowledge Discovery and Data
491 Mining*, pages 30–42. Springer, 2024.
- 492 [53] H. Kim, M. Kim, F. Berto, J. Kim, and J. Park. DevFormer: A symmetric transformer for
493 context-aware device placement. *International Conference on Machine Learning*, 2023.
- 494 [54] M. Kim, J. Park, and J. Kim. Learning collaborative policies to solve np-hard routing prob-
495 lems. In *Advances in Neural Information Processing Systems*, 2021.
- 496 [55] M. Kim, J. Park, and J. Park. Sym-NCO: Leveraging symmetry for neural combinatorial
497 optimization. *Advances in Neural Information Processing Systems*, 2022.
- 498 [56] M. Kim, S. Choi, J. Son, H. Kim, J. Park, and Y. Bengio. Ant colony sampling with
499 GFlowNets for combinatorial optimization. *arXiv preprint arXiv:2403.07041*, 2024.
- 500 [57] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint
501 arXiv:1412.6980*, 2014.
- 502 [58] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks.
503 In *International Conference on Learning Representations*, 2017.
- 504 [59] V. Konda and J. Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing
505 systems*, 12, 1999.
- 506 [60] W. Kool, H. Van Hoof, and M. Welling. Attention, learn to solve routing problems! *Interna-
507 tional Conference on Learning Representations*, 2019.

- 508 [61] W. Kool, H. Van Hoof, and M. Welling. Stochastic beams and where to find them: The
509 gumbel-top-k trick for sampling sequences without replacement. In *International Conference*
510 *on Machine Learning*, pages 3499–3508. PMLR, 2019.
- 511 [62] Y.-D. Kwon, J. Choo, B. Kim, I. Yoon, Y. Gwon, and S. Min. POMO: Policy optimization
512 with multiple optima for reinforcement learning. *Advances in Neural Information Processing*
513 *Systems*, 33:21188–21198, 2020.
- 514 [63] Y.-D. Kwon, J. Choo, I. Yoon, M. Park, D. Park, and Y. Gwon. Matrix encoding networks for
515 neural combinatorial optimization. *Advances in Neural Information Processing Systems*, 34:
516 5138–5149, 2021.
- 517 [64] G. Laporte and S. Martello. The selective travelling salesman problem. *Discrete applied*
518 *mathematics*, 26(2-3):193–207, 1990.
- 519 [65] E. Lawler, J. Lenstra, A. R. Kan, and D. Shmoys. The traveling salesman problem: A guided
520 tour of combinatorial optimization. *The Journal of the Operational Research Society*, 37(5):
521 535, 1986.
- 522 [66] G. Li, C. Xiong, A. Thabet, and B. Ghanem. Deepergcn: All you need to train deeper gens.
523 *arXiv preprint arXiv:2006.07739*, 2020.
- 524 [67] J. Li, L. Xin, Z. Cao, A. Lim, W. Song, and J. Zhang. Heterogeneous attentions for solv-
525 ing pickup and delivery problem via deep reinforcement learning. *IEEE Transactions on*
526 *Intelligent Transportation Systems*, 23(3):2306–2315, 2021.
- 527 [68] J. Li, Y. Ma, Z. Cao, Y. Wu, W. Song, J. Zhang, and Y. M. Chee. Learning feature embedding
528 refiner for solving vehicle routing problems. *IEEE Transactions on Neural Network and*
529 *Learning Systems*, 2023.
- 530 [69] S. Li, Y. Zhao, R. Varma, O. Salpekar, P. Noordhuis, T. Li, A. Paszke, J. Smith, B. Vaughan,
531 P. Damania, et al. Pytorch distributed: Experiences on accelerating data parallel training.
532 *arXiv preprint arXiv:2006.15704*, 2020.
- 533 [70] E. Liang, R. Liaw, R. Nishihara, P. Moritz, R. Fox, J. Gonzalez, K. Goldberg, and I. Sto-
534 ica. Ray rllib: A composable and scalable reinforcement learning library. *arXiv preprint*
535 *arXiv:1712.09381*, 85, 2017.
- 536 [71] J. T. Linderoth, A. Lodi, et al. Milp software. *Wiley encyclopedia of operations research and*
537 *management science*, 5:3239–3248, 2010.
- 538 [72] F. Liu, X. Lin, Q. Zhang, X. Tong, and M. Yuan. Multi-task learning for routing problem
539 with cross-problem zero-shot generalization. In *Proceedings of the 30th ACM SIGKDD Con-*
540 *ference on Knowledge Discovery and Data Mining*, 2024.
- 541 [73] F. Liu, X. Tong, M. Yuan, X. Lin, F. Luo, Z. Wang, Z. Lu, and Q. Zhang. Evolution of
542 heuristics: Towards efficient automatic algorithm design using large language model. In
543 *International Conference on Machine Learning*, 2024.
- 544 [74] R. Lotfi, A. Mostafaeipour, N. Mardani, and S. Mardani. Investigation of wind farm location
545 planning by considering budget constraints. *International Journal of Sustainable Energy*, 37
546 (8):799–817, 2018.
- 547 [75] F. Luo, X. Lin, F. Liu, Q. Zhang, and Z. Wang. Neural combinatorial optimization with heavy
548 decoder: Toward large scale generalization. *Advances in Neural Information Processing*
549 *Systems*, 36, 2024.
- 550 [76] F. Luo, X. Lin, Z. Wang, T. Xialiang, M. Yuan, and Q. Zhang. Self-improved learning for
551 scalable neural combinatorial optimization. *arXiv preprint arXiv:2403.19561*, 2024.

- 552 [77] L. Luttmann and L. Xie. Neural combinatorial optimization on heterogeneous graphs: An
553 application to the picker routing problem in mixed-shelves warehouses. In *Proceedings of*
554 *the International Conference on Automated Planning and Scheduling*, volume 34, pages 351–
555 359, 2024.
- 556 [78] Y. Ma, J. Li, Z. Cao, W. Song, L. Zhang, Z. Chen, and J. Tang. Learning to iteratively solve
557 routing problems with dual-aspect collaborative transformer. *Advances in Neural Information*
558 *Processing Systems*, 34, 2021.
- 559 [79] Y. Ma, J. Li, Z. Cao, W. Song, H. Guo, Y. Gong, and Y. M. Chee. Efficient neural neighbor-
560 hood search for pickup and delivery problems. *arXiv preprint arXiv:2204.11399*, 2022.
- 561 [80] Y. Ma, Z. Cao, and Y. M. Chee. Learning to search feasible and infeasible regions of routing
562 problems with flexible neural k-opt. *Advances in Neural Information Processing Systems*, 36,
563 2024.
- 564 [81] V. Makovychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller,
565 N. Rudin, A. Allshire, A. Handa, and G. State. I: High performance GPU-based physics
566 simulation for robot learning, 2021.
- 567 [82] S. Manchanda, S. Michel, D. Drakulic, and J.-M. Andreoli. On the generalization of neural
568 combinatorial optimization heuristics. In *Machine Learning and Knowledge Discovery in*
569 *Databases: European Conference, ECML PKDD 2022, Grenoble, France, September 19–23,*
570 *2022, Proceedings, Part V*, pages 426–442. Springer, 2023.
- 571 [83] V. Marianov, D. Serra, et al. Location problems in the public sector. *Facility location:*
572 *Applications and theory*, 1:119–150, 2002.
- 573 [84] Y. Min, Y. Bai, and C. P. Gomes. Unsupervised learning for solving the travelling salesman
574 problem. In *Neural Information Processing Systems*, 2023.
- 575 [85] V. Mnih, N. Heess, A. Graves, et al. Recurrent models of visual attention. *Advances in neural*
576 *information processing systems*, 27, 2014.
- 577 [86] V. Moens. TensorDict: your PyTorch universal data carrier, 2023. URL [https://github.](https://github.com/pytorch-labs/tensordict)
578 [com/pytorch-labs/tensordict](https://github.com/pytorch-labs/tensordict).
- 579 [87] S. A. Mulder and D. C. Wunsch II. Million city traveling salesman problem solution by divide
580 and conquer clustering with adaptive resonance neural networks. *Neural Networks*, 16(5-6):
581 827–832, 2003.
- 582 [88] A. T. Murray, K. Kim, J. W. Davis, R. Machiraju, and R. Parent. Coverage optimization to
583 support security monitoring. *Computers, Environment and Urban Systems*, 31(2):133–147,
584 2007.
- 585 [89] M. Nazari, A. Oroojlooy, L. Snyder, and M. Takác. Reinforcement learning for solving the
586 vehicle routing problem. *Advances in neural information processing systems*, 31, 2018.
- 587 [90] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agar-
588 wal, K. Slama, A. Ray, et al. Training language models to follow instructions with human
589 feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- 590 [91] M. Pagliardini, D. Paliotta, M. Jaggi, and F. Fleuret. Faster causal attention over large se-
591 quences through sparse flash attention. *arXiv preprint arXiv:2306.01160*, 2023.
- 592 [92] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,
593 N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learn-
594 ing library. *Advances in neural information processing systems*, 32, 2019.

- 595 [93] L. Perron and V. Furnon. OR-Tools, 2023. URL [https://developers.google.com/](https://developers.google.com/optimization/)
596 [optimization/](https://developers.google.com/optimization/).
- 597 [94] A. Prouvost, J. Dumouchelle, L. Scavuzzo, M. Gasse, D. Chételat, and A. Lodi. Ecole:
598 A gym-like library for machine learning in combinatorial optimization solvers. In *Learning*
599 *Meets Combinatorial Algorithms at NeurIPS2020*, 2020. URL [https://openreview.net/](https://openreview.net/forum?id=IVc9hqgibyB)
600 [forum?id=IVc9hqgibyB](https://openreview.net/forum?id=IVc9hqgibyB).
- 601 [95] C. PyVRP. PyVRP, 2023. URL <https://pyvrp.org/>.
- 602 [96] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann. Stable-
603 baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning*
604 *Research*, 22(268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.
- 605 [97] G. K. Rand. Sequencing and scheduling: An introduction to the mathematics of the job-
606 shop. *Journal of the Operational Research Society*, 33:862, 1982. URL [https://api.](https://api.semanticscholar.org/CorpusID:62592932)
607 [semanticscholar.org/CorpusID:62592932](https://api.semanticscholar.org/CorpusID:62592932).
- 608 [98] B. Romera-Paredes, M. Barekatin, A. Novikov, M. Balog, M. P. Kumar, E. Dupont, F. J.
609 Ruiz, J. S. Ellenberg, P. Wang, O. Fawzi, et al. Mathematical discoveries from program
610 search with large language models. *Nature*, 625(7995):468–475, 2024.
- 611 [99] M. W. Savelsbergh and M. Sol. The general pickup and delivery problem. *Transportation*
612 *science*, 29(1):17–29, 1995.
- 613 [100] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimiza-
614 tion algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 615 [101] W. Shan, Q. Yan, C. Chen, M. Zhang, B. Yao, and X. Fu. Optimization of competitive facility
616 location for chain stores. *Annals of Operations research*, 273:187–205, 2019.
- 617 [102] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously
618 large neural networks: The sparsely-gated mixture-of-experts layer. In *ICLR*, 2017.
- 619 [103] J. Son, M. Kim, S. Choi, and J. Park. Solving np-hard min-max routing problems as sequential
620 generation with equity context. *arXiv preprint arXiv:2306.02689*, 2023.
- 621 [104] J. Son, M. Kim, H. Kim, and J. Park. Meta-SAGE: Scale meta-learning scheduled adap-
622 tation with guided exploration for mitigating scale shift on combinatorial optimization. In
623 *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pages
624 32194–32210. PMLR, 2023.
- 625 [105] J. Song, Y. Yue, B. Dilkina, et al. A general large neighborhood search framework for solving
626 integer linear programs. *Advances in Neural Information Processing Systems*, 33:20012–
627 20023, 2020.
- 628 [106] W. Song, X. Chen, Q. Li, and Z. Cao. Flexible job-shop scheduling via graph neural network
629 and deep reinforcement learning. *IEEE Transactions on Industrial Informatics*, 19(2):1600–
630 1610, 2022.
- 631 [107] L. Sun, W. Huang, P. S. Yu, and W. Chen. Multi-round influence maximization. In *Pro-*
632 *ceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data*
633 *mining*, pages 2249–2258, 2018.
- 634 [108] Z. Sun and Y. Yang. Difusco: Graph-based diffusion solvers for combinatorial optimization.
635 *arXiv preprint arXiv:2302.08224*, 2023.
- 636 [109] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for rein-
637 forcement learning with function approximation. *Advances in neural information processing*
638 *systems*, 12, 1999.

- 639 [110] P. Tassel, M. Gebser, and K. Schekotihin. A reinforcement learning environment for job-shop
640 scheduling. *arXiv preprint arXiv:2104.03760*, 2021.
- 641 [111] D. Thyssens, T. Deredde, J. K. Falkner, and L. Schmidt-Thieme. Routing arena: A bench-
642 mark suite for neural routing solvers. *arXiv preprint arXiv:2310.04140*, 2023.
- 643 [112] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Lio, Y. Bengio, et al. Graph attention
644 networks. *stat*, 1050(20):10–48550, 2017.
- 645 [113] T. Vidal. Hybrid genetic search for the cvrp: Open-source implementation and swap* neigh-
646 borhood. *Computers & Operations Research*, 140:105643, 2022.
- 647 [114] O. Vinyals, M. Fortunato, and N. Jaitly. Pointer networks. In C. Cortes, N. Lawrence, D. Lee,
648 M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
649 volume 28, pages 2692–2700. Curran Associates, Inc., 2015. URL [https://proceedings.
650 neurips.cc/paper/2015/file/29921001f2f04bd3baee84a12e98098f-Paper.pdf](https://proceedings.neurips.cc/paper/2015/file/29921001f2f04bd3baee84a12e98098f-Paper.pdf).
- 651 [115] O. Vinyals, M. Fortunato, and N. Jaitly. Pointer networks. *Advances in neural information
652 processing systems*, 28, 2015.
- 653 [116] C. P. Wan, T. Li, and J. M. Wang. RLOR: A flexible framework of deep reinforcement
654 learning for operation research. *arXiv preprint arXiv:2303.13117*, 2023.
- 655 [117] R. Wang, L. Shen, Y. Chen, X. Yang, D. Tao, and J. Yan. Towards one-shot neural combi-
656 natorial solvers: Theoretical and empirical notes on the cardinality-constrained case. In *The
657 Eleventh International Conference on Learning Representations*, 2022.
- 658 [118] S. Wasserkrug, L. Boussioux, D. d. Hertog, F. Mirzazadeh, I. Birbil, J. Kurtz, and D. Maragno.
659 From large language models and optimization to decision optimization CoPilot: A research
660 manifesto. *arXiv preprint arXiv:2402.16269*, 2024.
- 661 [119] J. Weng, H. Chen, D. Yan, K. You, A. Duburcq, M. Zhang, Y. Su, H. Su, and J. Zhu. Tianshou:
662 A highly modularized deep reinforcement learning library. *Journal of Machine Learning
663 Research*, 23(267):1–6, 2022.
- 664 [120] Y. Wu, W. Song, Z. Cao, J. Zhang, and A. Lim. Learning improvement heuristics for solving
665 routing problems. *IEEE transactions on neural networks and learning systems*, 33(9):5057–
666 5069, 2021.
- 667 [121] Z. Xiao, D. Zhang, Y. Wu, L. Xu, Y. J. Wang, X. Han, X. Fu, T. Zhong, J. Zeng, M. Song,
668 and G. Chen. Chain-of-experts: When LLMs meet complex operations research problems. In
669 *International Conference on Learning Representations*, 2024.
- 670 [122] L. Xin, W. Song, Z. Cao, and J. Zhang. Generative adversarial training for neural combinato-
671 rial optimization models, 2022. URL <https://openreview.net/forum?id=9vsRT9mc7U>.
- 672 [123] O. Yadan. Hydra - a framework for elegantly configuring complex applications. Github,
673 2019. URL <https://github.com/facebookresearch/hydra>.
- 674 [124] C. Yang, X. Wang, Y. Lu, H. Liu, Q. V. Le, D. Zhou, and X. Chen. Large language models as
675 optimizers. In *International Conference on Learning Representations*, 2024.
- 676 [125] H. Ye, J. Wang, Z. Cao, H. Liang, and Y. Li. Deepaco: Neural-enhanced ant systems for
677 combinatorial optimization. *arXiv preprint arXiv:2309.14032*, 2023.
- 678 [126] H. Ye, J. Wang, Z. Cao, F. Berto, C. Hua, H. Kim, J. Park, and G. Song. Large language
679 models as hyper-heuristics for combinatorial optimization. *arXiv preprint arXiv:2402.01145*,
680 2024.

- 681 [127] H. Ye, J. Wang, H. Liang, Z. Cao, Y. Li, and F. Li. GLOP: Learning global partition and
682 local construction for solving large-scale routing problems in real-time. In *Proceedings of the*
683 *AAAI Conference on Artificial Intelligence*, volume 38, pages 20284–20292, 2024.
- 684 [128] C. Zhang, W. Song, Z. Cao, J. Zhang, P. S. Tan, and X. Chi. Learning to dispatch for job shop
685 scheduling via deep reinforcement learning. *Advances in Neural Information Processing*
686 *Systems*, 33:1621–1632, 2020.
- 687 [129] D. Zhang, H. Dai, N. Malkin, A. C. Courville, Y. Bengio, and L. Pan. Let the flows tell:
688 Solving graph combinatorial problems with gflownets. In A. Oh, T. Naumann, A. Globerson,
689 K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing*
690 *Systems*, volume 36, pages 11952–11969. Curran Associates, Inc., 2023.
- 691 [130] J. Zhou, Y. Wu, W. Song, Z. Cao, and J. Zhang. Towards omni-generalizable neural methods
692 for vehicle routing problems. In *International Conference on Machine Learning*, 2023.
- 693 [131] J. Zhou, Z. Cao, Y. Wu, W. Song, Y. Ma, J. Zhang, and C. Xu. MVMoE: Multi-task vehicle
694 routing solver with mixture-of-experts. In *International Conference on Machine Learning*,
695 2024.

696 **Checklist**

- 697 1. For all authors...
- 698 (a) Do the main claims made in the abstract and introduction accurately reflect the pa-
699 per’s contributions and scope? [Yes] Each claim has the corresponding contents in the
700 manuscript.
- 701 (b) Did you describe the limitations of your work? [Yes] See § 6.1.
- 702 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See § 7.
- 703 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
704 them? [Yes] We have read them and make sure that our paper conform to them.
- 705 2. If you are including theoretical results...
- 706 (a) Did you state the full set of assumptions of all theoretical results? [N/A] We do not
707 present theoretical results in this work.
- 708 (b) Did you include complete proofs of all theoretical results? [N/A] We do not present
709 theoretical results in this work.
- 710 3. If you ran experiments (e.g. for benchmarks)...
- 711 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
712 imental results (either in the supplemental material or as a URL)? [Yes] We made the
713 whole project open-sourced at <https://github.com/ai4co/rl4co>.
- 714 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
715 were chosen)? [Yes] See Appendix.
- 716 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
717 iments multiple times)? [Yes] We note that, as common practice in the field, we did
718 not report multiple runs for the main tables as algorithms can take more than one day
719 each to train. However, for experiments limited in the number of samples, such as for
720 the sample efficiency experiments and the mDPP benchmarking, we reported multiple
721 runs with different random seeds, where we demonstrated the robustness of different
722 runs to random seeds.
- 723 (d) Did you include the total amount of compute and the type of resources used (e.g., type
724 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix.
- 725 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 726 (a) If your work uses existing assets, did you cite the creators? [Yes] All the assets are
727 properly cited.
- 728 (b) Did you mention the license of the assets? [Yes] See Appendix.
- 729 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
730 All the new assets are available at <https://github.com/ai4co/rl4co>.
- 731 (d) Did you discuss whether and how consent was obtained from people whose data
732 you’re using/curating? [Yes] We discussed the licenses under which we obtained ac-
733 cess to the assets.
- 734 (e) Did you discuss whether the data you are using/curating contains personally identifi-
735 able information or offensive content? [N/A] The assets in this work do not involve
736 personally identifiable information or offensive content.
- 737 5. If you used crowdsourcing or conducted research with human subjects...
- 738 (a) Did you include the full text of instructions given to participants and screenshots, if
739 applicable? [N/A] To the best of our knowledge, this work does not involve crowd-
740 sourcing or human subjects.
- 741 (b) Did you describe any potential participant risks, with links to Institutional Review
742 Board (IRB) approvals, if applicable? [N/A] This work does not involve crowdsourc-
743 ing or human subjects.

744
745
746

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] This work does not involve crowdsourcing or human subjects.