

RL4CO: AN EXTENSIVE REINFORCEMENT LEARNING FOR COMBINATORIAL OPTIMIZATION BENCHMARK

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep reinforcement learning (RL) has recently shown significant benefits in solving combinatorial optimization (CO) problems, reducing reliance on domain expertise, and improving computational efficiency. However, the field lacks a unified benchmark for easy development and standardized comparison of algorithms across diverse CO problems. To fill this gap, we introduce RL4CO, a unified and extensive benchmark with in-depth library coverage of 23 state-of-the-art methods and more than 20 CO problems. Built on efficient software libraries and best practices in implementation, RL4CO features modularized implementation and flexible configuration of diverse RL algorithms, neural network architectures, inference techniques, and environments. RL4CO allows researchers to seamlessly navigate existing successes and develop their unique designs, facilitating the entire research process by decoupling science from heavy engineering. We also provide extensive benchmark studies to inspire new insights and future work. RL4CO has attracted numerous researchers in the community and is open-sourced.¹

1 INTRODUCTION

Combinatorial optimization (CO) focuses on finding optimal solutions for problems with discrete variables and has broad applications, including vehicle routing (Nazari et al., 2018; Kool et al., 2019a), scheduling (Zhang et al., 2020), and hardware device placement (Kim et al., 2023). Given that the combinatorial space expands exponentially and exhibits NP-hard characteristics, the operations research (OR) community has traditionally tackled these challenges through the development of mathematical programming algorithms (Gurobi Optimization, 2021) and handcrafted heuristics (Mart et al., 2018). Despite their success, these methods still face significant limitations: mathematical programming struggles with scaling, while handcrafted heuristics require significant domain-specific adjustments for different CO problems.

Recently, to address these limitations, neural combinatorial optimization (NCO) has emerged (Ben-gio et al., 2021b). It employs deep neural networks to automate the problem-solving process and significantly reduces computation demands and domain expertise requirements. Recent NCO works mainly leverage the reinforcement learning (RL) paradigm, making significant strides in improving exploration efficiency (Kwon et al., 2020; Kim et al., 2021), relaxing the needs of obtaining optimal solutions, and extending to various CO tasks (Zhang et al., 2020; Nazari et al., 2018; Kool et al., 2019a; Kim et al., 2023). Although supervised learning methods (Drakulic et al., 2023) are shown to be effective in NCO, they require a large number of high quality solutions, which is unrealistic for large instances or theoretically hard problems. Therefore, this work focuses on the RL paradigm.

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

¹Documentation: <https://anonymous.4open.science/w/rl4co-submission/>.
Code: <https://anonymous.4open.science/r/rl4co-submission/>.

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

Library	Environments #	Baselines [†] #	Hardware Acceleration	Availability	Modular Baselines	Open Community
ORL (Balaji et al., 2019)	3	1	×	×	×	×
OR-Gym (Hubbs et al., 2020)	9	-	×	✓	×	×
Graph-Env (Biagioni et al., 2022)	2	-	×	✓	×	×
RLOR (Wan et al., 2023)	2	2	×	✓	✓	×
RoutingArena (Thyssens et al., 2023)	1	8	✓	×	×	×
Jumanji (Bonnet et al., 2024)	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.

Contributions. To bridge this gap, we introduce RL4CO, the first comprehensive benchmark with multiple baselines, environments, and boilerplate from the literature, all implemented in a *modular, flexible, accelerated, and unified* manner. Our aim is to facilitate the entire research process for the NCO community with the following key contributions: 1) *Simplifying development* through modularizing 27 environments and 23 existing baseline models, allowing for flexible and automated combinations for effortless testing, switching, and achieving state-of-the-art performance; 2) *Enhancing the training and testing efficiency* through the customized unified pipeline tailored for the NCO community based on advanced libraries such as TorchRL (Bou et al., 2024), PyTorch Lightning (Falcon and The PyTorch Lightning team, 2019), Hydra (Yadan, 2019), and TensorDict (Moens, 2023); 3) *Standardizing evaluation* to ensure fair and comprehensive comparisons, enabling researchers to automatically test a broader range of problems from diverse distributions and gather valuable insights using our testbed. Overall, RL4CO eliminates the need for repetitive heavy engineering in the NCO community and fosters seamless future development by building on existing successes, enabling advanced innovation and progress in the field.

2 RELATED WORKS

Neural Combinatorial Optimization. Neural combinatorial optimization (NCO) utilizes machine learning techniques to automatically develop novel heuristics for solving NP-hard CO problems. We classify the majority of NCO research from the following perspectives: 1) *Learning Paradigms*: researchers have employed supervised learning (Vinyals et al., 2015; Hottung et al., 2020; Sun and Yang, 2023; Drakulic et al., 2023; Luo et al., 2024a) to approximate optimal solutions to CO instances. Further research leverages reinforcement learning (Bello et al., 2017; Nazari et al., 2018; Kool et al., 2019a; Kwon et al., 2020), and unsupervised learning (Min et al., 2023) to ease the difficulty of obtaining (near-)optimal solutions. 2) *Models*: various deep learning architectures such as recurrent neural networks (Vinyals et al., 2015; Chen and Tian, 2019; Li et al., 2023), graph neural networks (Joshi et al., 2019; Min et al., 2023), Transformers (Kool et al., 2019a; Kwon et al., 2020), diffusion models (Sun and Yang, 2023), and GFlowNets Zhang et al. (2023); Kim et al. (2024) have been employed. 3) *Problems*: NCO has demonstrated great success in various problems, including vehicle routing problems (VRPs) (e.g., traveling salesman problem (TSP) and capacitated VRP), scheduling problems (e.g., job shop scheduling problems (Zhang et al., 2020)), hardware device placement (Kim et al., 2023), and graph-based CO (Zhang et al., 2023). 4) *Heuristic Types*: in general, the learned heuristics can be categorized as *constructive* in an autoregressive (Kool et al., 2019a) or non-autoregressive (Joshi et al., 2019) way, and *improvement* heuristics that leverage traditional heuristics (Wu et al., 2021; Ma et al., 2024) and meta-heuristics (Song et al., 2020). We refer to Bengio et al. (2021b) for a comprehensive survey. In this paper, we focus on the reinforcement learning paradigm due to its effectiveness and flexibility. Notably, the proposed RL4CO is versatile to support most combinations of models, problems and heuristic types, making it an apt library and benchmark for future research in NCO.

Related Benchmark Libraries. Despite the variety of general-purpose RL software libraries (Brockman et al., 2016; Liang et al., 2017; Raffin et al., 2021; Weng et al., 2022) there is a lack of a unified and extensive benchmark for CO problems. Balaji et al. (2019) propose an RL benchmark for OR that comes only with a PPO baseline (Schulman et al., 2017). Also, Hubbs et al. (2020) and Biagioni et al. (2022) implement a collection of OR environments, but do not provide

any baselines to solve them. [Wan et al. \(2023\)](#) propose an RL for OR library that benchmarks the canonical TSP and CVRP environments using different configurations of the attention model ([Kool et al., 2019a](#)). Despite having a different focus, we also mention the work of [Prouvost et al. \(2020\)](#) here, who develop an API to use RL for controlling traditional MILP solvers ([Lindereth et al., 2010](#)), rather than using RL directly to learn solutions. Besides their narrow scope, a major downside of the above libraries is that they cannot be massively parallelized due to their reliance on the OpenAI Gym API ([Brockman et al., 2016](#)), which can only run on CPU. In contrast, RL4CO is based on TorchRL ([Bou et al., 2024](#)), a recent official PyTorch ([Paszke et al., 2019](#)) library for RL that enables hardware-accelerated execution of both environments and algorithms.

In contrast to the works above, Routing Arena ([Thyssens et al., 2023](#)) provides multiple neural and classical baselines, but benchmarks these only on the CVRP. The most related recent work is Jumanji ([Bonnet et al., 2024](#)), which provides a variety of CO environments written in JAX ([Bradbury et al., 2018](#)) that can be hardware-accelerated alongside an actor-critic baseline. While Jumanji is an RL environment suite, RL4CO is a full-stack library that integrates environments, policies, and RL algorithms under a unified framework. As such, baselines in RL4CO are modular and applicable to all suitable CO problems, whereas in Jumanji, policies are tailored to a specific environment.

3 RL4CO: TAXONOMY

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

Environments. Given a CO problem instance x , we formulate the solution-generating procedure as a Markov Decision Process (MDP) characterized by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$ as follows. **State** \mathcal{S} is the space of states that represent the given problem x and the current partial solution being updated in the MDP. **Action** \mathcal{A} is the action space, which includes all feasible actions a_t that can be taken at each step t . **State Transition** \mathcal{T} is the deterministic state transition function $s_{t+1} = \mathcal{T}(s_t, a_t)$ that 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 the immediate reward received after taking action a_t in state s_t . Finally, $\gamma \in [0, 1]$ is a discount factor that determines the importance of future rewards. Since the state transition is deterministic, we represent the solution for a problem x as a sequence of T actions $\mathbf{a} = (a_0, \dots, a_{T-1})$. Then, the total return $\sum_{t=0}^{T-1} \mathcal{R}(s_t, a_t)$ translates to the negative cost function of the CO problem.

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

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

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

A NAR policy encodes a problem x into a heuristic $\mathcal{H} = f(x) \in \mathbb{R}_+^N$, where N is the number of possible assignments across all decision variables. Each number in \mathcal{H} represents an (unnormalized) probability of a particular assignment. To obtain a solution \mathbf{a} from \mathcal{H} , one can sample a sequence of assignments from \mathcal{H} while dynamically masking infeasible assignments to meet problem-specific constraints. It can also guide a search process, e.g., Ant Colony Optimization ([Dorigo and Stützle, 2019](#); [Ye et al., 2023](#); [Kim et al., 2024](#)), or be incorporated into hybrid frameworks ([Ye et al., 2024b](#)). Here, the heuristic helps identify promising transitions and improve the efficiency of finding an optimal or near-optimal solution.

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

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

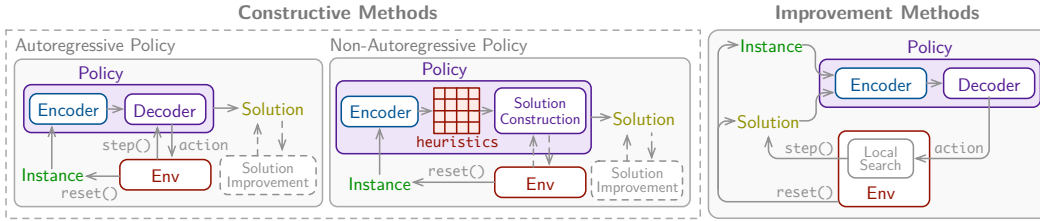


Figure 1: Overview of different types of policies and their modularization in RL4CO.

where \mathbf{a}^k is the k -th updated solution and K is the budget in terms of the number of improvements. This process allows continuous refinement over time to improve the quality of the solution.

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

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

where θ is the set of parameters of π and $P(\mathbf{x})$ is the distribution of problem instances. Eq. (3) can be solved using algorithms such as variations of REINFORCE (Sutton et al., 1999), Advantage Actor-Critic (A2C) methods (Konda and Tsitsiklis, 1999), or Proximal Policy Optimization (PPO) (Schulman et al., 2017). These algorithms are employed to train the policy network π , by transforming the maximization problem in Eq. (3) into a minimization problem involving a loss function, which is then optimized using gradient descent algorithms. For instance, the REINFORCE loss function gradient is given by:

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

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

4 RL4CO: LIBRARY STRUCTURE

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

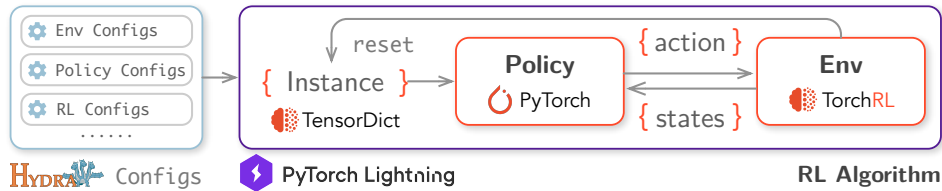


Figure 2: Overview of the RL4CO pipeline: from configurations to training a policy.

4.1 ENVIRONMENTS

Environment Interface. Environments in RL4CO fully specify the CO problems and their logic in a stateless manner. That is, all static and dynamic information about the environment like the current state s_t , actions a_t , and rewards r_t are passed to and retrieved from the environment’s `reset` and `step` functions through a `TensorDict` (Moens, 2023). This not only enables modular in-

216 interactions between environments and policies, but also allows for seamless integration with various
 217 components without the need for environment-specific adaptations. Further, a key advantage com-
 218 pared to other libraries is that environments in RL4CO can take batches of instances and process
 219 them in parallel on a GPU. Instances \mathbf{x} are provided through a modular `generator` to the envi-
 220 ronment, and different generators can be used to generate different data distributions, facilitating the
 221 generalization of trained policies.

222 RL4CO’s environments are based on the `RL4COEnvBase` class that extends the `EnvBase` of
 223 TorchRL (Bou et al., 2024) with additional features and efficiency improvements. For instance,
 224 RL4CO’s `step` method brings a decrease of up to 50% in latency and halves the memory impact
 225 by keeping only required transitions in the stateless `TensorDict`. Additionally, our environment
 226 API contains several functions, such as `render` and `check_solution_validity`, helping to
 227 analyze generated solutions, `select_start_nodes` for multi-start methods like POMO (Kwon
 228 et al., 2020) and `local_search` for iterative solution improvement.

229 **Problems.** We include benchmarking for 27 environments divided into four areas. 1) *Routing*:
 230 Traveling Salesman Problem (TSP) (Lawler et al., 1986), Capacitated Vehicle Routing Problem
 231 (CVRP) (Bodin, 1983), Orienteering Problem (OP) (Laporte and Martello, 1990; Chao et al., 1996),
 232 Prize Collecting TSP (PCTSP) (Balas, 1989), Pickup and Delivery Problem (PDP) (Kalantari et al.,
 233 1985; Savelsbergh and Sol, 1995) and Multi-Task VRP (MTVRP), which includes 16 problem vari-
 234 ants, namely the basic VRPTW, OVRP, VRPB, VRPL and VRPs with the respective constraint combi-
 235 nations (Liu et al., 2024a; Zhou et al., 2024; Berto et al., 2024); 2) *Scheduling*: Flexible Job Shop
 236 Scheduling Problem (FJSSP) (Brandimarte, 1993), Job Shop Scheduling Problem (JSSP) (Rand,
 237 1982) and Flexible Flow Shop Problem (FFSP); 3) *Electronic Design Automation*: multiple De-
 238 cap Placement Problem (mDPP) (Kim et al., 2023); 4) *Graph*: Facility Location Problem (FLP)
 239 (Drezner and Hamacher, 2004) and Max Cover Problem (MCP) (Khuller et al., 1999). A detailed
 240 description of the environment implementations for these problems can be found in Appendix B.

241 4.2 POLICIES

242 Policies in RL4CO are subclasses of PyTorch’s `nn.Module` and contain the encoding-decoding
 243 logic and neural network parameters θ . Drawing on our taxonomy in § 3, RL4CO pro-
 244 vides different metaclasses like `AutoregressivePolicy`, `NonAutoregressivePolicy`,
 245 or `ImprovementPolicy` that the different policies in the RL4CO “zoo” can inherit from.
 246 RL4CO modularize components to process environment specific features into the embedding space
 247 via parametrized functions. First, *node embeddings* $\phi_n : \mathbb{R}^{N \times m_n} \rightarrow \mathbb{R}^{N \times h}$ transform m_n raw
 248 features for the N nodes of problem \mathbf{x} from the feature space to the embedding space h . Further,
 249 *edge embeddings* $\phi_e : \mathbb{R}^{E \times m_e} \rightarrow \mathbb{R}^{E \times h}$ transform m_e edge features of instance \mathbf{x} from the fea-
 250 ture space to the embedding space h , where E is the number of edges. Lastly, *context embeddings*
 251 $\phi_c : \mathbb{R}^{m_c} \rightarrow \mathbb{R}^h$ capture contextual information not related to a specific node or edge by transform-
 252 ing m_c context features from the current decoding step s_t from the feature space to the embedding
 253 space h . Overall, Fig. 3 illustrates a generic constructive AR policy in RL4CO, where the feature
 254 embeddings are applied similarly to other types of policies. Feature projections can be automatically
 255 selected by RL4CO at runtime by simply passing the environment to the policy. Additionally, we
 256 allow for granular control of any higher-level policy component, such as encoders and decoders.

258 4.3 RL ALGORITHMS

259 RL algorithms in RL4CO are used to learn the parameters θ of the `Policy` by interact-
 260 ing with the `Environment` and its problem instances. The parent class of algorithms is the
 261 `RL4COLitModule`, inheriting from PyTorch Lightning’s `pl.LightningModule` (Falcon and
 262 The PyTorch Lightning team, 2019). This allows for granular support of various methods includ-
 263 ing the `[train, val, test]_step`, automatic logging with several logging services such as
 264 Wandb via `log_metrics`, automatic optimizer configuration via `configure_optimizers`
 265 and several useful callbacks for RL methods such as `on_train_epoch_end`. RL algorithms
 266 are additionally attached to an `RL4COTrainer`, a wrapper we made with additional optimizations
 267 around `pl.Trainer`. This module seamlessly supports features of modern training pipelines, in-
 268 cluding logging, checkpoint management, mixed-precision training, various hardware acceleration
 269 supports (e.g., CPU, GPU, TPU, and Apple Silicon), and multi-device distribution (Li et al., 2020b).

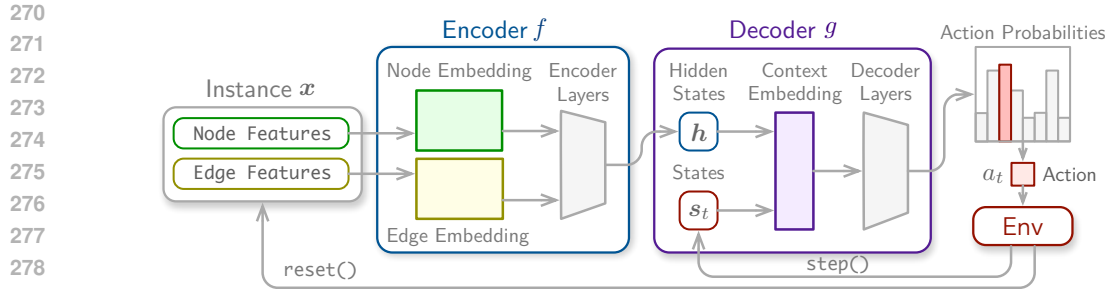


Figure 3: Overview of modularized RL4CO policies. Any component such as the encoder/decoder structure and feature embeddings can be replaced and thus the model is adaptable to various new environments.

For instance, using mixed-precision training significantly decreases training time without sacrificing much convergence and enables us to leverage recent routines, e.g., FlashAttention (Dao et al., 2022; Dao, 2023), which we investigate in Appendix E.7.2.

4.4 UTILITIES

Configuration Management. Optional, but useful, RL4CO utilizes Hydra (Yadan, 2019), a framework that enables hierarchical configuration management. This not only makes it easy to define complex configurations and manage lots of experiments with different (hyperparameter) settings, but also facilitates reproducibility of experiments by freezing experimental setups in configuration files. We outline the process of configuration management with Hydra in Appendix D.3.1.

Decoding Schemes. Decoding schemes define the logic of translating from the model’s unnormalized log-probabilities to actions. Specifically, they handle the transformation from the unnormalized log-probabilities to probabilities $P(\mathcal{A})$ by masking infeasible actions and optionally applying tanh clipping (Bello et al., 2017) prior to softmax normalization. Subsequently, different decoding strategies can be employed to determine the action based on the probability distribution: 1) *Greedy*, which selects the action with the highest probability; 2) *Sampling*, which samples from the masked probability distribution of the policy, where different sampling strategies like softmax temperature scaling τ , top- k sampling (Kool et al., 2019b), and top- p (or Nucleus) sampling (Holtzman et al., 2019) can be used; 3) *Multistart*, which enforces diverse starting actions as demonstrated in POMO (Kwon et al., 2020), such as starting from different cities in the TSP; 4) *Augmentation*, which applies transformations to instances, such as random rotations in Euclidean problems (Kim et al., 2022), to create an augmented set of problems. We describe these strategies in detail in Appendix D.4.

Documentation, Tutorials, and Testing. We release extensive documentation² to make RL4CO as accessible as possible for both newcomers and experts. Furthermore, several tutorials and examples are also available under the `examples/` folder of our publicly available code³. We thoroughly test our library via continuous integration on multiple Python versions and operating systems. The following code snippet shows minimalistic code that can train a model in a few lines:

```

311 from rl4co.envs.routing import TSPEnv, TSPGenerator
312 from rl4co.models import AttentionModelPolicy, POMO
313 from rl4co.utils import RL4COTrainer
314 # Instantiate generator and environment
315 generator = TSPGenerator(num_loc=50, loc_distribution="uniform")
316 env = TSPEnv(generator)
317 # Create policy and RL model
318 policy = AttentionModelPolicy(env_name=env.name, num_encoder_layers=6)
319 model = POMO(env, policy, batch_size=64)
320 # Instantiate Trainer and fit
321 trainer = RL4COTrainer(epochs=10, accelerator="gpu", precision="fp16")
322 trainer.fit(model)

```

²<https://anonymous.4open.science/w/rl4co-submission/>

³<https://anonymous.4open.science/r/rl4co-submission/>

4.5 BASELINES ZOO

RL4CO entails a collection of 23 neural solvers for CO from the literature, which are implemented in a modular and flexible way using our policy metaclasses introduced in § 4.2. We organize these baselines into five categories: (non-)autoregressive constructive methods, improvement strategies, general-purpose RL algorithms and transductive RL methods. Table 2 lists all baselines currently available in RL4CO and we refer the reader to Appendix C for further implementation details.

Table 2: RL4CO baselines.

Category	Methods
Constructive AR Methods	AM (Kool et al., 2019a), Ptr-Net (Vinyals et al., 2015), POMO (Kwon et al., 2020), HAM (Li et al., 2021), SymNCO (Kim et al., 2022), PolyNet (Hottung et al., 2024), MVMoE (Zhou et al., 2024), L2D (Zhang et al., 2020), HGNN (Song et al., 2022), MatNet (Kwon et al., 2021), MTPO (Liu et al., 2024a), DevFormer (Kim et al., 2023)
Constructive NAR Methods	DeepACO (Ye et al., 2023), GFACS (Kim et al., 2024), GLOP (Ye et al., 2024b)
Improvement Methods	DACT (Ma et al., 2021), N2S (Ma et al., 2022), NeuOpt (Ma et al., 2024)
RL Algorithms	REINFORCE (Sutton et al., 1999), Advantage Actor-Critic (A2C) (Konda and Tsitsiklis, 1999), Proximal Policy Optimization (PPO) (Schulman et al., 2017)
Transductive RL Methods	Active search (Bello et al., 2017), Efficient active search (Hottung et al., 2022)

5 BENCHMARKING STUDY

We perform several benchmarking studies with our unified RL4CO library, with experimental setup and benchmarking details in Appendix D. Due to the extent of our benchmark, we report highlights of the results in the following section and refer the reader to Appendix E for additional experiments.

5.1 FLEXIBILITY AND MODULARITY

Changing policy components. The integration of many state-of-the-art methods in RL4CO from the NCO field in a modular framework makes it easy to implement and improve upon state-of-the-art neural solvers for complex CO problems with only a few lines of code.⁴ We demonstrate this in Table 3 for the FJSSP by gradually replacing or adding elements to the original SotA policy (Song et al., 2022). First, replacing the HGNN encoder with the more expressive MatNet encoder (Kwon et al., 2021) already improves the average makespan by around 7%. Further, replacing the MLP decoder with the pointer mechanism from the AM model (Kool et al., 2019a) reduces the optimality gap to roughly one-third of that observed in the policy proposed by Song et al. (2022), even when using a greedy decoding strategy.

Table 3: Solutions obtained with RL4CO for the FJSSP with different model configurations.

Encoder / Decoder		FJSSP	
		10 × 5	20 × 5
HGNN + MLP (g.) (Song et al., 2022)	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%

5.2 CONSTRUCTIVE POLICIES

Mind Your Baseline. In on-policy RL, which is often employed in RL4CO due to fast reward function evaluations, several different REINFORCE baselines have been proposed to improve the performance. We benchmark several RL algorithms training constructive policies for routing problems of node size 50, whose underlying architecture is based on the encoder-decoder Attention Model (Kool et al., 2019a) and whose main difference lies in how the REINFORCE baseline is calculated (we additionally train the AM with PPO as further reference). For a fair comparison,

⁴The different model configurations shown here can be obtained by simply changing the Hydra configuration file like the one shown in Appendix D.

we run all baselines in controlled settings with the same number of optimization steps and report results in Table 4. The performance of A2C is generally inferior to other baselines. This can be explained by the inherent challenge of estimating the value of a problem instance x based on the sparse reward, which is only observed after solving an entire instance in routing problems. We found similar trends regarding actor-critic methods as A2C and PPO in the EDA mDPP problem (Kim et al., 2023), which we report in Appendix E.4. Namely, a greedy rollout baseline (Kool et al., 2019a) can do better than value-based methods due to the challenging task of instance value estimation.

Table 4: Optimality gaps 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

Interestingly, while POMO (Kwon et al., 2020), which takes as a baseline the average reward over all routes forcing each starting node to be different, may work well as baselines for problems in which near-optimal solutions can be constructed from any node (e.g., TSP), this may not be true for other problems such as the Orienteering Problem (OP): the reason is that in OP only a subset of nodes should be selected in

an optimal solution, while several states will be discarded. Hence, forcing the policy to select all of them makes up for a poor baseline. We remark that while SymNCO (whose shared baseline involves symmetric rotations and flips) (Kim et al., 2022) may perform well in Euclidean problems, it is not applicable in non-Euclidean CO, including asymmetric routing problems and scheduling.

Decoding Schemes. The solution quality of different solvers often shows significant improvements in performance with different decoding schemes. We evaluate the pre-trained solver with different strategies and settings as shown in Fig. 4.

Generalization. Using RL4CO, we can easily evaluate the generalization performance of existing baselines by employing supported environments that incorporate various VRP variant tasks and instance distributions (termed MTPOMO and MDPOMO, respectively). Empirical results on CVRPLib, shown in Table 5, reveal that training on different tasks significantly enhances generalization performance. This finding underscores the necessity of building foundation models across diverse CO domains.

Large-Scale Instances. We evaluate large-scale CVRP instances of thousands of nodes, with more visualizations and scaling in Appendix E.1.6. The last row of Table 6 illustrates the performance of the hybrid NAR/AR GLOP (Ye et al., 2024b), while others refer to reproduced results from Ye et al. (2024b). Our implementation in RL4CO improves the performance in not only speed but also solution quality.

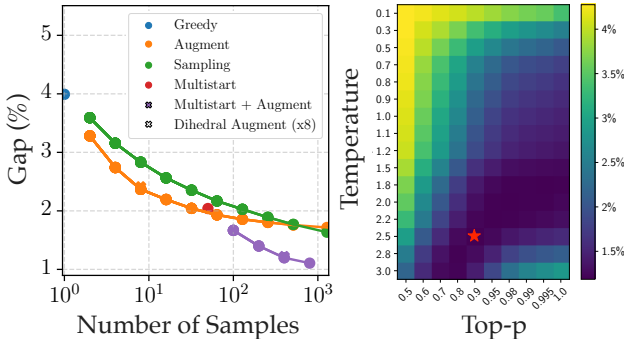


Figure 4: Study of decoding schemes using POMO on CVRP50. [Left]: Pareto front of decoding schemes by the number of samples; [Right]: performance of sampling with different temperatures τ and p values for top- p sampling.

Table 5: Results on CVRPLIB with models trained on $N = 50$. Greedy multi-start decoding is used.

	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 6: Performance on large-scale CVRP instances with thousands of nodes.

	CVRP1K		CVRP2K		CVRP7K	
	Obj.	Time	Obj.	Time	Obj.	Time
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

5.3 COMBINING CONSTRUCTION AND IMPROVEMENT: BEST OF BOTH WORLDS?

While constructive policies can build solutions in seconds, their performance is often limited, even with advanced decoding schemes such as sampling or augmentations. On the other hand, improvement methods are more suitable for larger computing budgets. We benchmark models on TSP with 50 nodes: the AR constructive method POMO (Kwon et al., 2020) and the improvement methods DACT (Ma et al., 2021) and NeuOpt (Ma et al., 2024). In the original implementation, DACT and NeuOpt started from a solution constructed randomly. To further demonstrate the flexibility of RL4CO, we show that bootstrapping improvement methods with constructive ones enhance convergence speed. Fig. 5 shows that bootstrapping with a pre-trained POMO policy significantly enhances the convergence speed. To further investigate the performance, we report the Primal Integral (PI) (Berthold, 2013; Vidal, 2022; Thyssens et al., 2023), which evaluates the evolution of solution quality over time. Improvement methods alone, such as DACT and NeuOpt, achieve 2.99 and 2.26 respectively, while sampling from POMO achieves 0.08. This shows that the “area under the curve” can be better even if the final solution is worse for constructive methods. Bootstrapping with POMO then improves DACT and NeuOpt to 0.08 and 0.04 respectively, showing the benefits of modularity and hybridization of different components.

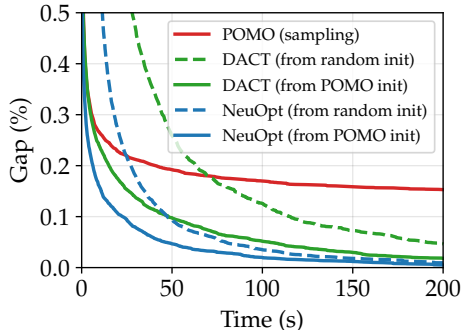


Figure 5: Bootstrapping improvement with constructive methods.

6 DISCUSSION

Limitations and Future Directions While RL4CO is an efficient and modular library specialized in CO problems, it might not be suitable for any other task due to a number of area-specific optimizations, and we do not expect it to seamlessly integrate with, for instance, OpenAI Gym wrappers without some modifications. Another limitation of the library is its scope so far, namely RL. Eventually, creating a new library to support supervised methods as a comprehensive “AI4CO” codebase could benefit the whole NCO community. We additionally identify in Foundation Models for CO and related scalable architectures a promising area of future research to overcome generalization issues across tasks and distributions, for which we provided some early clues in Appendix E.8.

Long-term Plans RL4CO is an active library that has already garnered much attention from the community, with over 400 stars on GitHub. We thank contributors in the community who have helped us build RL4CO. Our long-term plan is to become the go-to RL for CO benchmark library. For this purpose, we created a community Slack workspace (link available upon acceptance) that has attracted more than 200 researchers. We are committed to helping resolve issues and questions from the community and actively engaged in discussion. It is our hope that our work will ultimately benefit the NCO field with new ideas and collaborations. More available in Appendix A.

7 CONCLUSION

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

REFERENCES

- 486
487
488 L. Accorsi, A. Lodi, and D. Vigo. Guidelines for the computational testing of machine learning
489 approaches to vehicle routing problems. *Operations Research Letters*, 50(2):229–234, 2022.
- 490
491 A. AhmadiTeshnizi, W. Gao, and M. Udell. OptiMUS: Scalable optimization modeling with (MI)LP
492 solvers and large language models. In *International Conference on Machine Learning*, 2024.
- 493
494 K. Ali, W. Alsalih, and H. Hassanein. Set-cover approximation algorithms for load-aware readers
495 placement in RFID networks. In *2011 IEEE international conference on communications (ICC)*,
pages 1–6. IEEE, 2011.
- 496
497 D. Applegate, R. Bixby, V. Chvátal, and W. Cook. Concorde TSP solver, 2023. URL <https://www.math.uwaterloo.ca/tsp/concorde/index.html>.
- 498
499 B. Balaji, J. Bell-Masterson, E. Bilgin, A. Damianou, P. M. Garcia, A. Jain, R. Luo, A. Maggiar,
500 B. Narayanaswamy, and C. Ye. Orl: Reinforcement learning benchmarks for online stochastic
501 optimization problems. *arXiv preprint arXiv:1911.10641*, 2019.
- 502
503 E. Balas. The prize collecting traveling salesman problem. *Networks*, 19(6):621–636, 1989.
- 504
505 A. Bdeir, J. K. Falkner, and L. Schmidt-Thieme. Attention, filling in the gaps for generalization in
506 routing problems. In *Joint European Conference on Machine Learning and Knowledge Discovery
507 in Databases*, pages 505–520. Springer, 2022.
- 508
509 I. Bello, H. Pham, Q. V. Le, M. Norouzi, and S. Bengio. Neural combinatorial optimization with
reinforcement learning, 2017.
- 510
511 E. Bengio, M. Jain, M. Korablyov, D. Precup, and Y. Bengio. Flow network based generative mod-
512 els for non-iterative diverse candidate generation. *Advances in Neural Information Processing
513 Systems*, 34:27381–27394, 2021a.
- 514
515 Y. Bengio, A. Lodi, and A. Prouvost. Machine learning for combinatorial optimization: a method-
ological tour d’horizon. *European Journal of Operational Research*, 290(2):405–421, 2021b.
- 516
517 Y. Bengio, S. Lahlou, T. Deleu, E. J. Hu, M. Tiwari, and E. Bengio. Gflownet foundations. *Journal
518 of Machine Learning Research*, 24(210):1–55, 2023.
- 519
520 T. Berthold. Measuring the impact of primal heuristics. *Operations Research Letters*, 41(6):611–
521 614, 2013.
- 522
523 F. Berto, C. Hua, N. G. Zepeda, A. Hottung, N. Wouda, L. Lan, K. Tierney, and J. Park.
RouteFinder: Towards foundation models for vehicle routing problems. *Arxiv*, 2024. URL
524 <https://github.com/ai4co/routefinder>.
- 525
526 K. Bestuzheva, M. Besançon, W.-K. Chen, A. Chmiela, T. Donkiewicz, J. van Doornmalen, L. Eifler,
527 O. Gaul, G. Gamrath, A. Gleixner, et al. The SCIP optimization suite 8.0. *arXiv 2112.08872*,
528 2021.
- 529
530 M. Bettini, A. Prorok, and V. Moens. Benchmarl: Benchmarking multi-agent reinforcement learn-
ing. *arXiv preprint arXiv:2312.01472*, 2023.
- 531
532 J. Bi, Y. Ma, J. Wang, Z. Cao, J. Chen, Y. Sun, and Y. M. Chee. Learning generalizable models for
533 vehicle routing problems via knowledge distillation. *Advances in Neural Information Processing
534 Systems*, 35:31226–31238, 2022.
- 535
536 D. Biagioni, C. E. Tripp, S. Clark, D. Duplyakin, J. Law, and P. C. S. John. graphenv: a python
537 library for reinforcement learning on graph search spaces. *Journal of Open Source Software*, 7
(77):4621, 2022.
- 538
539 L. Bodin. Routing and scheduling of vehicles and crews. *Computer & Operations Research*, 10(2):
69–211, 1983.

- 540 C. Bonnet, D. Luo, D. Byrne, S. Surana, S. Abramowitz, P. Duckworth, V. Coyette, L. I. Midgley,
541 E. Tegegn, T. Kalloniatis, O. Mahjoub, M. Macfarlane, A. P. Smit, N. Grinsztajn, R. Boige, C. N.
542 Waters, M. A. Mimouni, U. A. M. Sob, R. de Kock, S. Singh, D. Furelos-Blanco, V. Le, A. Pre-
543 torius, and A. Laterre. Jumanji: a diverse suite of scalable reinforcement learning environments
544 in jax. In *International Conference on Learning Representations*, 2024.
- 545 A. Bou, M. Bettini, S. Dittert, V. Kumar, S. Sodhani, X. Yang, G. D. Fabritiis, and V. Moens.
546 TorchRL: A data-driven decision-making library for pytorch. In *International conference on*
547 *learning representations*, 2024.
- 549 J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke,
550 J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transformations of
551 Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- 552 P. Brandimarte. Routing and scheduling in a flexible job shop by tabu search. *Annals of Operations*
553 *research*, 41(3):157–183, 1993.
- 555 G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai
556 gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 557 S. Brody, U. Alon, and E. Yahav. How attentive are graph attention networks? In *International*
558 *Conference on Learning Representations*, 2019.
- 560 F. Bu, H. Jo, S. Y. Lee, S. Ahn, and K. Shin. Tackling prevalent conditions in unsupervised com-
561 binatorial optimization: Cardinality, minimum, covering, and more. In *International Conference*
562 *on Machine Learning*, 2024.
- 563 B. Çatay. Ant colony optimization and its application to the vehicle routing problem with pickups
564 and deliveries. In *Natural intelligence for scheduling, Planning and packing problems*, pages
565 219–244. Springer, 2009.
- 567 F. Chalumeau, S. Surana, C. Bonnet, N. Grinsztajn, A. Pretorius, A. Laterre, and T. Barrett. Com-
568 binatorial optimization with policy adaptation using latent space search. *Advances in Neural*
569 *Information Processing Systems*, 36, 2024.
- 570 I.-M. Chao, B. L. Golden, and E. A. Wasil. A fast and effective heuristic for the orienteering
571 problem. *European journal of operational research*, 88(3):475–489, 1996.
- 573 J. Chen, Z. Zhang, Z. Cao, Y. Wu, Y. Ma, T. Ye, and J. Wang. Neural multi-objective combinatorial
574 optimization with diversity enhancement. *Advances in Neural Information Processing Systems*,
575 36, 2024.
- 576 X. Chen and Y. Tian. Learning to perform local rewriting for combinatorial optimization. In *Ad-*
577 *vances in Neural Information Processing Systems*, 2019.
- 579 T. Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. *arXiv*
580 *preprint arXiv:2307.08691*, 2023.
- 581 T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré. Flashattention: Fast and memory-efficient exact
582 attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–
583 16359, 2022.
- 585 M. Dorigo and T. Stützle. *Ant colony optimization: overview and recent advances*. Springer, 2019.
- 586 D. Drakulic, S. Michel, F. Mai, A. Sors, and J.-M. Andreoli. BQ-NCO: Bisimulation quotienting
587 for generalizable neural combinatorial optimization. *Advances in Neural Information Processing*
588 *Systems*, 2023.
- 590 Z. Drezner and H. W. Hamacher. *Facility location: applications and theory*. Springer Science &
591 Business Media, 2004.
- 592 W. Falcon and The PyTorch Lightning team. PyTorch Lightning, 3 2019. URL <https://github.com/Lightning-AI/lightning>.

- 594 J. K. Falkner and L. Schmidt-Thieme. Learning to solve vehicle routing problems with time windows
595 through joint attention. *arXiv preprint arXiv:2006.09100*, 2020.
- 596
- 597 M. Fischetti, J. J. S. Gonzalez, and P. Toth. Solving the orienteering problem through branch-and-
598 cut. *INFORMS Journal on Computing*, 10(2):133–148, 1998.
- 599 N. Grinsztajn, D. Furelos-Blanco, S. Surana, C. Bonnet, and T. Barrett. Winner takes it all: Train-
600 ing performant rl populations for combinatorial optimization. *Advances in Neural Information*
601 *Processing Systems*, 36:48485–48509, 2023.
- 602
- 603 L. Gurobi Optimization. Gurobi optimizer reference manual, 2021. URL [http://www.gurobi.](http://www.gurobi.com)
604 [com](http://www.gurobi.com).
- 605 W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. *Advances*
606 *in neural information processing systems*, 30, 2017.
- 607
- 608 K. Helsgaun. An extension of the Lin-Kernighan-Helsgaun TSP solver for constrained traveling
609 salesman and vehicle routing problems. *Roskilde: Roskilde University*, 12 2017. doi: 10.13140/
610 RG.2.2.25569.40807.
- 611 S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780,
612 1997.
- 613
- 614 A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi. The curious case of neural text degeneration.
615 *arXiv preprint arXiv:1904.09751*, 2019.
- 616 A. Hottung, B. Bhandari, and K. Tierney. Learning a latent search space for routing problems using
617 variational autoencoders. In *International Conference on Learning Representations*, 2020.
- 618
- 619 A. Hottung, Y.-D. Kwon, and K. Tierney. Efficient active search for combinatorial optimization
620 problems. *International conference on learning representations*, 2022.
- 621
- 622 A. Hottung, M. Mahajan, and K. Tierney. PolyNet: Learning diverse solution strategies for neural
623 combinatorial optimization. *arXiv preprint arXiv:2402.14048*, 2024.
- 624
- 625 Y. Hou, H. Ye, Y. Zhang, S. Xu, and G. Song. Routeplacer: An end-to-end routability-aware placer
626 with graph neural network. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge*
627 *Discovery and Data Mining*, 2024.
- 628
- 629 C. D. Hubbs, H. D. Perez, O. Sarwar, N. V. Sahinidis, I. E. Grossmann, and J. M. Wassick.
OR-Gym: A reinforcement learning library for operations research problems. *arXiv preprint*
arXiv:2008.06319, 2020.
- 630
- 631 J. Hwang, J. S. Pak, D. Yoon, H. Lee, J. Jeong, Y. Heo, and I. Kim. Enhancing on-die pdn for optimal
632 use of package pdn with decoupling capacitor. In *2021 IEEE 71st Electronic Components and*
633 *Technology Conference (ECTC)*, pages 1825–1830, 2021. doi: 10.1109/ECTC32696.2021.00288.
- 634
- 635 Z. Iklassov, Y. Du, F. Akimov, and M. Takac. Self-guiding exploration for combinatorial problems.
arXiv preprint arXiv:2405.17950, 2024.
- 636
- 637 S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing
638 internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr,
2015.
- 639
- 640 R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton. Adaptive mixtures of local experts.
641 *Neural computation*, 3(1):79–87, 1991.
- 642
- 643 A. D. Jesus, A. Liefoghe, B. Derbel, and L. Paquete. Algorithm selection of anytime algorithms. In
Proceedings of the 2020 genetic and evolutionary computation conference, pages 850–858, 2020.
- 644
- 645 Y. Jiang, Y. Wu, Z. Cao, and J. Zhang. Learning to solve routing problems via distributionally robust
646 optimization. In *36th AAAI Conference on Artificial Intelligence*, 2022.
- 647
- M. I. Jordan and R. A. Jacobs. Hierarchical mixtures of experts and the em algorithm. *Neural*
computation, 6(2):181–214, 1994.

- 648 C. K. Joshi, T. Laurent, and X. Bresson. An efficient graph convolutional network technique for the
649 travelling salesman problem. *arXiv preprint arXiv:1906.01227*, 2019.
650
- 651 J. Juang, L. Zhang, Z. Kiguradze, B. Pu, S. Jin, and C. Hwang. A modified genetic algorithm for the
652 selection of decoupling capacitors in pdn design. In *2021 IEEE International Joint EMC/SI/PI and*
653 *EMC Europe Symposium*, pages 712–717, 2021. doi: 10.1109/EMC/SI/PI/EMCEurope52599.
654 2021.9559292.
- 655 B. Kalantari, A. V. Hill, and S. R. Arora. An algorithm for the traveling salesman problem with
656 pickup and delivery customers. *European Journal of Operational Research*, 22(3):377–386, 1985.
657
- 658 E. Khalil, H. Dai, Y. Zhang, B. Dilkina, and L. Song. Learning combinatorial optimization algo-
659 rithms over graphs. *Advances in neural information processing systems*, 30, 2017.
660
- 661 S. Khuller, A. Moss, and J. S. Naor. The budgeted maximum coverage problem. *Information*
662 *processing letters*, 70(1):39–45, 1999.
- 663 D. Kikuta, H. Ikeuchi, K. Tajiri, and Y. Nakano. RouteExplainer: An explanation framework for
664 vehicle routing problem. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*,
665 pages 30–42. Springer, 2024.
666
- 667 H. Kim, M. Kim, F. Berto, J. Kim, and J. Park. DevFormer: A symmetric transformer for context-
668 aware device placement. *International Conference on Machine Learning*, 2023.
- 669 M. Kim, J. Park, and J. Kim. Learning collaborative policies to solve NP-hard routing problems. In
670 *Advances in Neural Information Processing Systems*, 2021.
671
- 672 M. Kim, J. Park, and J. Park. Sym-NCO: Leveraging symmetricity for neural combinatorial opti-
673 mization. *Advances in Neural Information Processing Systems*, 2022.
674
- 675 M. Kim, S. Choi, J. Son, H. Kim, J. Park, and Y. Bengio. Ant colony sampling with GFlowNets for
676 combinatorial optimization. *arXiv preprint arXiv:2403.07041*, 2024.
- 677 D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint*
678 *arXiv:1412.6980*, 2014.
679
- 680 T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In
681 *International Conference on Learning Representations*, 2017.
- 682 V. Konda and J. Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing*
683 *systems*, 12, 1999.
684
- 685 W. Kool, H. Van Hoof, and M. Welling. Attention, learn to solve routing problems! *International*
686 *Conference on Learning Representations*, 2019a.
687
- 688 W. Kool, H. Van Hoof, and M. Welling. Stochastic beams and where to find them: The gumbel-
689 top-k trick for sampling sequences without replacement. In *International Conference on Machine*
690 *Learning*, pages 3499–3508. PMLR, 2019b.
- 691 Y.-D. Kwon, J. Choo, B. Kim, I. Yoon, Y. Gwon, and S. Min. POMO: Policy optimization with
692 multiple optima for reinforcement learning. *Advances in Neural Information Processing Systems*,
693 33:21188–21198, 2020.
694
- 695 Y.-D. Kwon, J. Choo, I. Yoon, M. Park, D. Park, and Y. Gwon. Matrix encoding networks for neural
696 combinatorial optimization. *Advances in Neural Information Processing Systems*, 34:5138–5149,
697 2021.
- 698 G. Laporte and S. Martello. The selective travelling salesman problem. *Discrete applied mathemat-*
699 *ics*, 26(2-3):193–207, 1990.
700
- 701 E. Lawler, J. Lenstra, A. R. Kan, and D. Shmoys. The traveling salesman problem: A guided tour of
combinatorial optimization. *The Journal of the Operational Research Society*, 37(5):535, 1986.

- 702 F. Li, B. Golden, and E. Wasil. The open vehicle routing problem: Algorithms, large-scale test
703 problems, and computational results. *Computers & Operations Research*, 34(10):2918–2930,
704 2007. ISSN 0305-0548. doi: <https://doi.org/10.1016/j.cor.2005.11.018>.
705
- 706 G. Li, C. Xiong, A. Thabet, and B. Ghanem. Deepergcn: All you need to train deeper gcns. *arXiv*
707 *preprint arXiv:2006.07739*, 2020a.
- 708 J. Li, L. Xin, Z. Cao, A. Lim, W. Song, and J. Zhang. Heterogeneous attentions for solving pickup
709 and delivery problem via deep reinforcement learning. *IEEE Transactions on Intelligent Trans-*
710 *portation Systems*, 23(3):2306–2315, 2021.
- 711 J. Li, Y. Ma, Z. Cao, Y. Wu, W. Song, J. Zhang, and Y. M. Chee. Learning feature embedding
712 refiner for solving vehicle routing problems. *IEEE Transactions on Neural Network and Learning*
713 *Systems*, 2023.
- 714 S. Li, Y. Zhao, R. Varma, O. Salpekar, P. Noordhuis, T. Li, A. Paszke, J. Smith, B. Vaughan, P. Da-
715 mania, et al. Pytorch distributed: Experiences on accelerating data parallel training. *arXiv preprint*
716 *arXiv:2006.15704*, 2020b.
- 717 E. Liang, R. Liaw, R. Nishihara, P. Moritz, R. Fox, J. Gonzalez, K. Goldberg, and I. Stoica. Ray rllib:
718 A composable and scalable reinforcement learning library. *arXiv preprint arXiv:1712.09381*, 85,
719 2017.
- 720 I. Lima, E. Uchoa, D. Pecin, A. Pessoa, M. Poggi, T. Vidal, A. Subramanian, R. W. D. Oliveira,
721 and E. Queiroga. CVRPLIB: Capacitated vehicle routing problem library, 2014. URL <http://vrp.galgos.inf.puc-rio.br/index.php/en/>. Last checked on October 6, 2024.
722
- 723 X. Lin, Z. Yang, and Q. Zhang. Pareto set learning for neural multi-objective combinatorial opti-
724 mization. *arXiv preprint arXiv:2203.15386*, 2022.
725
- 726 J. T. Linderoth, A. Lodi, et al. Milp software. *Wiley encyclopedia of operations research and*
727 *management science*, 5:3239–3248, 2010.
728
- 729 F. Liu, X. Lin, Q. Zhang, X. Tong, and M. Yuan. Multi-task learning for routing problem with
730 cross-problem zero-shot generalization. In *Proceedings of the 30th ACM SIGKDD Conference on*
731 *Knowledge Discovery and Data Mining*, 2024a.
- 732 F. Liu, X. Tong, M. Yuan, X. Lin, F. Luo, Z. Wang, Z. Lu, and Q. Zhang. Evolution of heuris-
733 tics: Towards efficient automatic algorithm design using large language model. In *International*
734 *Conference on Machine Learning*, 2024b.
- 735 S. Liu, C. Chen, X. Qu, K. Tang, and Y.-S. Ong. Large language models as evolutionary optimizers.
736 *arXiv preprint arXiv:2310.19046*, 2023.
737
- 738 R. Lotfi, A. Mostafaeipour, N. Mardani, and S. Mardani. Investigation of wind farm location plan-
739 ning by considering budget constraints. *International Journal of Sustainable Energy*, 37(8):799–
740 817, 2018.
- 741 F. Luo, X. Lin, F. Liu, Q. Zhang, and Z. Wang. Neural combinatorial optimization with heavy
742 decoder: Toward large scale generalization. *Advances in Neural Information Processing Systems*,
743 36, 2024a.
- 744 F. Luo, X. Lin, Z. Wang, T. Xialiang, M. Yuan, and Q. Zhang. Self-improved learning for scalable
745 neural combinatorial optimization. *arXiv preprint arXiv:2403.19561*, 2024b.
- 746 L. Luttmann and L. Xie. Neural combinatorial optimization on heterogeneous graphs: An applica-
747 tion to the picker routing problem in mixed-shelves warehouses. In *Proceedings of the Interna-*
748 *tional Conference on Automated Planning and Scheduling*, volume 34, pages 351–359, 2024.
- 749 Y. Ma, J. Li, Z. Cao, W. Song, L. Zhang, Z. Chen, and J. Tang. Learning to iteratively solve routing
750 problems with dual-aspect collaborative transformer. *Advances in Neural Information Processing*
751 *Systems*, 34, 2021.

- 756 Y. Ma, J. Li, Z. Cao, W. Song, H. Guo, Y. Gong, and Y. M. Chee. Efficient neural neighborhood
757 search for pickup and delivery problems. *arXiv preprint arXiv:2204.11399*, 2022.
- 758
- 759 Y. Ma, Z. Cao, and Y. M. Chee. Learning to search feasible and infeasible regions of routing
760 problems with flexible neural k-opt. *Advances in Neural Information Processing Systems*, 36,
761 2024.
- 762 S. Manchanda, S. Michel, D. Drakulic, and J.-M. Andreoli. On the generalization of neural combi-
763 natorial optimization heuristics. In *Machine Learning and Knowledge Discovery in Databases:
764 European Conference, ECML PKDD 2022, Grenoble, France, September 19–23, 2022, Proceed-
765 ings, Part V*, pages 426–442. Springer, 2023.
- 766 V. Marianov, D. Serra, et al. Location problems in the public sector. *Facility location: Applications
767 and theory*, 1:119–150, 2002.
- 768
- 769 R. Mart, P. M. Pardalos, and M. G. Resende. *Handbook of heuristics*. Springer Publishing Company,
770 Incorporated, 2018.
- 771 Y. Min, Y. Bai, and C. P. Gomes. Unsupervised learning for solving the travelling salesman problem.
772 In *Neural Information Processing Systems*, 2023.
- 773
- 774 V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Ried-
775 miller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement
776 learning. *nature*, 518(7540):529–533, 2015.
- 777 V. Moens. TensorDict: your PyTorch universal data carrier, 2023. URL [https://github.
778 com/pytorch-labs/tensordict](https://github.com/pytorch-labs/tensordict).
- 779
- 780 A. T. Murray, K. Kim, J. W. Davis, R. Machiraju, and R. Parent. Coverage optimization to support
781 security monitoring. *Computers, Environment and Urban Systems*, 31(2):133–147, 2007.
- 782 M. Nazari, A. Oroojlooy, L. Snyder, and M. Takác. Reinforcement learning for solving the vehicle
783 routing problem. *Advances in neural information processing systems*, 31, 2018.
- 784
- 785 M. Pagliardini, D. Paliotta, M. Jaggi, and F. Fleuret. Faster causal attention over large sequences
786 through sparse flash attention. *arXiv preprint arXiv:2306.01160*, 2023.
- 787 H. Park, H. Kim, H. Kim, J. Park, S. Choi, J. Kim, K. Son, H. Suh, T. Kim, J. Ahn, et al. Versatile
788 genetic algorithm-bayesian optimization (ga-bo) bi-level optimization for decoupling capacitor
789 placement. In *2023 IEEE 32nd Conference on Electrical Performance of Electronic Packaging
790 and Systems (EPEPS)*, pages 1–3. IEEE, 2023a.
- 791
- 792 J. Park, C. Kwon, and J. Park. Learn to solve the min-max multiple traveling salesmen problem with
793 reinforcement learning. In *Proceedings of the 2023 International Conference on Autonomous
794 Agents and Multiagent Systems*, pages 878–886, 2023b.
- 795 A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein,
796 L. Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. *Advances
797 in neural information processing systems*, 32, 2019.
- 798
- 799 L. Perron and V. Furnon. OR-Tools, 2023. URL [https://developers.google.com/
800 optimization/](https://developers.google.com/optimization/).
- 801 J. Pirnay and D. G. Grimm. Self-improvement for neural combinatorial optimization: Sample with-
802 out replacement, but improvement. *arXiv preprint arXiv:2403.15180*, 2024.
- 803
- 804 A. Prouvost, J. Dumouchelle, L. Scavuzzo, M. Gasse, D. Chételat, and A. Lodi. Ecole: A gym-like
805 library for machine learning in combinatorial optimization solvers. In *Learning Meets Combinatorial
806 Algorithms at NeurIPS2020*, 2020. URL [https://openreview.net/forum?id=
807 IVc9hqqibyB](https://openreview.net/forum?id=IVc9hqqibyB).
- 808 A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann. Stable-baselines3:
809 Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22
(268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.

- 810 G. K. Rand. Sequencing and scheduling: An introduction to the mathematics of the job-
811 shop. *Journal of the Operational Research Society*, 33:862, 1982. URL [https://api.
812 semanticscholar.org/CorpusID:62592932](https://api.semanticscholar.org/CorpusID:62592932).
- 813 G. Reinelt. TspLib—a traveling salesman problem library. *ORSA journal on computing*, 3(4):376–
814 384, 1991.
- 815
- 816 B. Romera-Paredes, M. Barekatin, A. Novikov, M. Balog, M. P. Kumar, E. Dupont, F. J. Ruiz, J. S.
817 Ellenberg, P. Wang, O. Fawzi, et al. Mathematical discoveries from program search with large
818 language models. *Nature*, 625(7995):468–475, 2024.
- 819 M. W. Savelsbergh and M. Sol. The general pickup and delivery problem. *Transportation science*,
820 29(1):17–29, 1995.
- 821
- 822 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
823 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 824 W. Shan, Q. Yan, C. Chen, M. Zhang, B. Yao, and X. Fu. Optimization of competitive facility
825 location for chain stores. *Annals of Operations research*, 273:187–205, 2019.
- 826
- 827 N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously
828 large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference
829 on Learning Representations*, 2017.
- 830 J. Son, M. Kim, H. Kim, and J. Park. Meta-SAGE: Scale meta-learning scheduled adaptation with
831 guided exploration for mitigating scale shift on combinatorial optimization. In *Proceedings of the
832 40th International Conference on Machine Learning*, volume 202, pages 32194–32210. PMLR,
833 2023.
- 834 J. Son, M. Kim, S. Choi, H. Kim, and J. Park. Equity-Transformer: Solving NP-hard min-max rout-
835 ing problems as sequential generation with equity context. In *Proceedings of the AAAI Conference
836 on Artificial Intelligence*, volume 38, pages 20265–20273, 2024.
- 837
- 838 J. Song, Y. Yue, B. Dilkina, et al. A general large neighborhood search framework for solving integer
839 linear programs. *Advances in Neural Information Processing Systems*, 33:20012–20023, 2020.
- 840 W. Song, X. Chen, Q. Li, and Z. Cao. Flexible job-shop scheduling via graph neural network and
841 deep reinforcement learning. *IEEE Transactions on Industrial Informatics*, 19(2):1600–1610,
842 2022.
- 843
- 844 L. Sun, W. Huang, P. S. Yu, and W. Chen. Multi-round influence maximization. In *Proceedings of
845 the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages
846 2249–2258, 2018.
- 847 Z. Sun and Y. Yang. DIFUSCO: Graph-based diffusion solvers for combinatorial optimization. In
848 *Advances in Neural Information Processing Systems*, volume 36, pages 3706–3731, 2023.
- 849 R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement
850 learning with function approximation. *Advances in neural information processing systems*, 12,
851 1999.
- 852
- 853 E. Taillard. Benchmarks for basic scheduling problems. *European journal of operational research*,
854 64(2):278–285, 1993.
- 855 H. Tang, F. Berto, Z. Ma, C. Hua, K. Ahn, and J. Park. Himap: Learning heuristics-informed policies
856 for large-scale multi-agent pathfinding. *arXiv preprint arXiv:2402.15546*, 2024a.
- 857
- 858 H. Tang, F. Berto, and J. Park. Ensembling prioritized hybrid policies for multi-agent pathfinding.
859 *arXiv preprint arXiv:2403.07559*, 2024b.
- 860 P. Tassel, M. Gebser, and K. Schekotihin. A reinforcement learning environment for job-shop
861 scheduling. *arXiv preprint arXiv:2104.03760*, 2021.
- 862
- 863 D. Thyssens, T. Deredde, J. K. Falkner, and L. Schmidt-Thieme. Routing arena: A benchmark
suite for neural routing solvers. *arXiv preprint arXiv:2310.04140*, 2023.

- 864 H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal,
865 E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint*
866 *arXiv:2302.13971*, 2023.
- 867 D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast
868 stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- 870 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin.
871 Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- 872 P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Lio, Y. Bengio, et al. Graph attention
873 networks. *stat*, 1050(20):10–48550, 2017.
- 875 T. Vidal. Hybrid genetic search for the CVRP: Open-source implementation and SWAP* neighbor-
876 hood. *Computers & Operations Research*, 140:105643, 2022.
- 877 O. Vinyals, M. Fortunato, and N. Jaitly. Pointer networks. In *Advances in Neural Information*
878 *Processing Systems*, volume 28, pages 2692–2700, 2015.
- 880 C. P. Wan, T. Li, and J. M. Wang. RLOR: A flexible framework of deep reinforcement learning for
881 operation research. *arXiv preprint arXiv:2303.13117*, 2023.
- 883 R. Wang, L. Shen, Y. Chen, X. Yang, D. Tao, and J. Yan. Towards one-shot neural combinatorial
884 solvers: Theoretical and empirical notes on the cardinality-constrained case. In *The Eleventh*
885 *International Conference on Learning Representations*, 2022.
- 886 S. Wasserkrug, L. Boussioux, D. d. Hertog, F. Mirzazadeh, I. Birbil, J. Kurtz, and D. Maragno. From
887 large language models and optimization to decision optimization CoPilot: A research manifesto.
888 *arXiv preprint arXiv:2402.16269*, 2024.
- 889 J. Weng, H. Chen, D. Yan, K. You, A. Duburcq, M. Zhang, Y. Su, H. Su, and J. Zhu. Tianshou: A
890 highly modularized deep reinforcement learning library. *Journal of Machine Learning Research*,
891 23(267):1–6, 2022.
- 893 N. A. Wouda, L. Lan, and W. Kool. PyVRP: A high-performance vrp solver package. *INFORMS*
894 *Journal on Computing*, 2024.
- 896 Y. Wu, W. Song, Z. Cao, J. Zhang, and A. Lim. Learning improvement heuristics for solving routing
897 problems. *IEEE transactions on neural networks and learning systems*, 33(9):5057–5069, 2021.
- 898 Z. Xiao, D. Zhang, Y. Wu, L. Xu, Y. J. Wang, X. Han, X. Fu, T. Zhong, J. Zeng, M. Song, and
899 G. Chen. Chain-of-experts: When LLMs meet complex operations research problems. In *Inter-*
900 *national Conference on Learning Representations*, 2024.
- 902 L. Xin, W. Song, Z. Cao, and J. Zhang. Generative adversarial training for neural combinatorial
903 optimization models, 2022. URL <https://openreview.net/forum?id=9vsRT9mc7U>.
- 904 O. Yadan. Hydra - a framework for elegantly configuring complex applications. Github, 2019. URL
905 <https://github.com/facebookresearch/hydra>.
- 907 C. Yang, X. Wang, Y. Lu, H. Liu, Q. V. Le, D. Zhou, and X. Chen. Large language models as
908 optimizers. In *International Conference on Learning Representations*, 2024.
- 909 H. Ye, J. Wang, Z. Cao, H. Liang, and Y. Li. DeepACO: Neural-enhanced ant systems for combina-
910 torial optimization. *arXiv preprint arXiv:2309.14032*, 2023.
- 912 H. Ye, J. Wang, Z. Cao, F. Berto, C. Hua, H. Kim, J. Park, and G. Song. Reevo: Large language mod-
913 els as hyper-heuristics with reflective evolution. In *Advances in Neural Information Processing*
914 *Systems*, 2024a. <https://github.com/ai4co/reevo>.
- 915 H. Ye, J. Wang, H. Liang, Z. Cao, Y. Li, and F. Li. GLOP: Learning global partition and local
916 construction for solving large-scale routing problems in real-time. In *Proceedings of the AAAI*
917 *Conference on Artificial Intelligence*, volume 38, pages 20284–20292, 2024b.

918 C. Zhang, W. Song, Z. Cao, J. Zhang, P. S. Tan, and X. Chi. Learning to dispatch for job shop
919 scheduling via deep reinforcement learning. *Advances in Neural Information Processing Systems*,
920 33:1621–1632, 2020.

921 D. Zhang, H. Dai, N. Malkin, A. C. Courville, Y. Bengio, and L. Pan. Let the flows tell: Solving
922 graph combinatorial problems with gflownets. In A. Oh, T. Naumann, A. Globerson, K. Saenko,
923 M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36,
924 pages 11952–11969. Curran Associates, Inc., 2023.

925 Z. Zheng, S. Yao, Z. Wang, X. Tong, M. Yuan, and K. Tang. Dpn: Decoupling partition and naviga-
926 tion for neural solvers of min-max vehicle routing problems. *arXiv preprint arXiv:2405.17272*,
927 2024.

928 J. Zhou, Y. Wu, W. Song, Z. Cao, and J. Zhang. Towards omni-generalizable neural methods for
929 vehicle routing problems. In *International Conference on Machine Learning*, 2023.

930 J. Zhou, Z. Cao, Y. Wu, W. Song, Y. Ma, J. Zhang, and C. Xu. MVMoE: Multi-task vehicle routing
931 solver with mixture-of-experts. In *International Conference on Machine Learning*, 2024.

932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

RL4CO: AN EXTENSIVE REINFORCEMENT LEARNING FOR COMBINATORIAL OPTIMIZATION BENCHMARK

Supplementary Material

Table of Contents

972		
973		
974		
975		
976		
977		
978		
979		
980		
981		
982	A	RL4CO: Vision and Software 21
983	A.1	Why Choosing the RL4CO Library? 21
984	A.2	On the Choice of the Software 21
985	A.3	Licenses 22
986	B	Environments 22
987	B.1	Routing 23
988	B.1.1	Traveling Salesman Problem (TSP) 23
989	B.1.2	Capacitated Vehicle Routing Problem (CVRP) 23
990	B.1.3	Orienteering Problem (OP) 24
991	B.1.4	Prize Collecting TSP (PCTSP) 24
992	B.1.5	Pickup and Delivery Problem (PDP) 24
993	B.1.6	Multi-Task VRP (MTVRP) 25
994	B.2	Scheduling 27
995	B.2.1	Job Shop Scheduling Problem (JSSP) 27
996	B.2.2	Flexible Job Shop Scheduling Problem (FJSSP) 27
997	B.2.3	Flexible Flow Shop Problem (FFSP) 28
998	B.3	Electronic Design Automation 28
999	B.3.1	Decap Placement Problem (DPP) 28
1000	B.3.2	Multi-Port Decap Placement Problem (mDPP) 29
1001	B.4	Graph 29
1002	B.4.1	Facility Location Problem (FLP) 29
1003	B.4.2	Maximum Coverage Problem (MCP) 30
1004	B.5	Additional Environments and Beyond 30
1005	C	Baselines 31
1006	C.1	General-purpose RL Algorithms 31
1007	C.1.1	REINFORCE 31
1008	C.1.2	Advantage Actor-Critic (A2C) 31
1009	C.1.3	Proximal Policy Optimization (PPO) 32
1010	C.2	Constructive Autoregressive (AR) 32
1011	C.2.1	Attention Model (AM) 32
1012	C.2.2	Ptr-Net 34
1013	C.2.3	POMO 34
1014	C.2.4	SymNCO 34
1015	C.2.5	PolyNet 35
1016	C.2.6	HAM 35
1017	C.2.7	MTPOMO 35
1018	C.2.8	MVMoE 35
1019	C.2.9	L2D 36
1020	C.2.10	HGNN 36
1021	C.2.11	MatNet 36
1022	C.2.12	DevFormer 37
1023	C.3	Constructive Non-Autoregressive (NAR) 37
1024	C.3.1	DeepACO 37
1025	C.3.2	GFACTS 38
	C.3.3	GLOP 38
	C.4	Improvement methods 39
	C.4.1	DACT 39
	C.4.2	N2S 39
	C.4.3	NeuOpt 39
	C.5	Active Search Methods 40
	C.5.1	Active Search (AS) 40

1026	C.5.2 Efficient Active Search (EAS)	40
1027		
1028	D Benchmarking Setup	40
1029	D.1 Metrics	40
1030	D.1.1 Gap to BKS	40
1031	D.1.2 Primal Integral	40
1032	D.1.3 Runtime Measurement	41
1033	D.2 Hardware & Software	41
1034	D.2.1 Hardware	41
1035	D.2.2 Software	42
1036	D.3 Hyperparameters	42
1037	D.3.1 Common Hyperparameters	42
1038	D.3.2 Changing Policy Components	42
1039	D.3.3 Mind Your Baseline	42
1040	D.3.4 Generalization: Cross-Task and Cross-Distribution	44
1041	D.3.5 Large-Scale Instances	44
1042	D.3.6 Combining Construction and Improvement	44
1043	D.4 Decoding Schemes	45
1044	D.4.1 Augmentations	45
1045	D.4.2 Sampling	45
1046		
1047	E Additional Experiments	47
1048	E.1 Mind your Baseline: Further Insights	47
1049	E.1.1 Main In-distribution Results	47
1050	E.1.2 Decoding Schemes Comparison	48
1051	E.1.3 Sample Efficiency	48
1052	E.1.4 Out-of-distribution	49
1053	E.1.5 Search Methods	50
1054	E.1.6 Additional Large-scale Results	52
1055	E.2 Learning Heuristics for Ant Colony Optimization	53
1056	E.2.1 Experiment Settings	53
1057	E.2.2 Results	53
1058	E.3 Learning to Schedule	53
1059	E.3.1 JSSP	54
1060	E.3.2 FJSSP	55
1061	E.3.3 FFSP	56
1062	E.3.4 Dense and Episodic Rewards	57
1063	E.4 Electronic Design Automation: Learning to Place Decaps	57
1064	E.4.1 Main Results	58
1065	E.4.2 Generalization to Different Number of Components	58
1066	E.5 Learning to Improve	58
1067	E.5.1 Main results	59
1068	E.5.2 Discussion	59
1069	E.6 Graph Problems: Facility Location Problem (FLP) and Maximum Coverage Problem (MCP)	60
1070	E.6.1 Experimental settings	60
1071	E.6.2 Benchmark Results	60
1072	E.6.3 Out-of-distribution	62
1073	E.7 Efficient Software Routines	65
1074	E.7.1 Mixed-Precision Training	65
1075	E.7.2 FlashAttention	67
1076	E.7.3 Efficient Memory Handling in Environments	67
1077	E.8 Towards Foundation Models	68
1078	E.8.1 Experimental Setting	68
1079	E.8.2 Empirical Results	69
	E.8.3 Discussion	69
	E.9 Generalization of Training on Multiple Distributions and Multiple Tasks	70

A RL4CO: VISION AND SOFTWARE

A.1 WHY CHOOSING THE RL4CO LIBRARY?

RL4CO is a *unified* and *extensive* benchmark for the RL-for-CO research domain, designed to be accessible and valuable to researchers and practitioners across all levels of expertise.



Figure 6: RL4CO benchmark logo.

Availability and Future Support RL4CO can be installed through PyPI. We adhere to continuous integration, deployment, and testing to ensure reproducibility and accessibility.

Open License We adopt the open MIT license for all content contained in RL4CO. We ascribe to the principles of *libre software*⁵. Most reimplementations are from original authors and are re-licensed under the MIT license. Data and baseline-specific licenses are reported in [Appendix A.3](#).



Figure 7: Unofficial - but widely used - open MIT license logo.

Open Community Through our journey, we started the AI4CO community⁶, which is a non-profit, cross-institution, inclusive, and open research community. AI4CO originally started out as a Slack channel for discussing the RL4CO but evolved into a broader-visioned and inclusive space to communicate with other researchers about general NCO. We warmly invite all interested people to join us.



Figure 8: AI4CO community logo.

A.2 ON THE CHOICE OF THE SOFTWARE

During the development of RL4CO, we wanted to make it as simple as possible to integrate reproducible and standardized code adhering to the latest guidelines. As a main template for our codebase, we use Lightning-Hydra-Template⁷ which we believe is a solid starting point for reproducible deep learning. We further discuss framework choices below.

PyTorch PyTorch (Paszke et al., 2019) is a popular open-source deep-learning framework that has gained significant traction in the research community. We chose PyTorch as the primary framework for RL4CO due to its intuitive API, dynamic computational graphs, strong community support, and seamless integration with the Python ecosystem. These features make PyTorch well-suited for

⁵<https://www.gnu.org/philosophy/free-sw.en.html>

⁶Community Github link available upon acceptance

⁷<https://github.com/ashleve/lightning-hydra-template>

rapid prototyping and experimentation, which are essential in research settings. Moreover, most of the existing research in NCO has been implemented. It is currently being implemented using PyTorch, making it not only easier to build upon and compare with previous work but also easier for newcomers and experienced researchers.

TorchRL and TensorDict One of the software hindrances in RL is the bottleneck between CPU and GPU communication, majorly due to CPU-based operating environments. For this reason, we did not opt for OpenAI Gym (Brockman et al., 2016) since, although it includes some level of parallelization, this does not happen on GPU and would thus greatly hinder performance. Kool et al. (2019a) creates *ad-hoc* environments in PyTorch to handle batched data efficiently. However, it could be cumbersome to integrate into standardized routines that include `step` and `reset` functions. As we searched for a better alternative, we found that TorchRL library (Bou et al., 2024), an official PyTorch project that allows for efficient batched implementations on (multiple) GPUs as well as functions akin to OpenAI Gym. We also employ the TensorDict (Bou et al., 2024) to handle tensors efficiently on multiple keys (i.e. in CVRP, we can directly operate transforms on multiple keys as locations, capacities, and more). This makes our environments compatible with the models in TorchRL, which we believe could further spread interest in the CO area.

PyTorch Lightning PyTorch Lightning (Falcon and The PyTorch Lightning team, 2019) is a useful tool for abstracting away the boilerplate code, allowing researchers and practitioners to focus more on the core ideas and innovations. It features a standardized training loop and an extensive set of pre-built components, including automated checkpointing, distributed training, and logging. PyTorch Lightning accelerates development time and facilitates scalability. We employ PyTorch Lightning in RL4CO to integrate with the PyTorch ecosystem - which includes TorchRL- enabling us to leverage the rich set of tools and libraries available.

Hydra Hydra (Yadan, 2019) is a powerful open-source framework for managing complex configurations in machine-learning models and other software. Hydra facilitates creating hierarchical configurations, making it easy to manage even very large and intricate configurations. Moreover, it integrates with command-line interfaces, allowing the execution of different configurations directly from the command line, thereby enhancing reproducibility. We found Hydra to be effective when dealing with multiple experiments since configurations are saved both locally, as `yaml` files, and can be uploaded to monitoring software as Wandb⁸ (or to any of the monitoring software supported by PyTorch Lightning).

A.3 LICENSES

Table 7: Reference code licenses and links.

Type	Asset	License	Link
Library	PyTorch Paszke et al. (2019)	BSD-3 License	link
	PyTorch Lightning Falcon and The PyTorch Lightning team (2019)	Apache-2.0 License	link
	TorchRL+TensorDict Bou et al. (2024)	MIT License	link
	Hydra Yadan (2019)	MIT License	link
Dataset	TSPLIB Reinelt (1991)	Available for any non-commercial use	link
	CVRLib Lima et al. (2014)	Available for any non-commercial use	link
	DPP PDNs (Park et al., 2023a)	Apache-2.0	link
Solver	PyVRP Wouda et al. (2024)	MIT	link
	LKH3 Helsgaun (2017)	Available for any non-commercial use	link
	OR-Tools Perron and Furnon (2023)	Apache 2.0 License	link

We summarize the license of software that we employ in RL4CO in a non-exhaustive list in Table 7. Original environments and models from the authors are acknowledged through their respective citations, with several links available in the library. RL4CO is licensed under the MIT license.

B ENVIRONMENTS

This section provides an overview of the list of environments we experimented with at the time of writing. We organize environments by categories, which, at the time of writing, are:

⁸<https://wandb.ai/>

1. **Routing (B.1)**
2. **Scheduling (B.2)**
3. **Electronic Design Automation (B.3)**
4. **Graph (B.4)**

B.1 ROUTING

Routing problems are perhaps the most known class of CO problems. They are problems of great practical importance, not only for logistics, where they are more commonly framed, but also for industry, engineering, science, and medicine. The typical objective of routing problems is to minimize the total length of the paths needed to visit some (or all) the nodes in a graph. In the following section, we present each of these variants with details of their implementations.

Common instance generation details Following the standard protocol of NCO for routing, we randomly sample node coordinates from the 2D unit square (i.e., $[0, 1]^2$). To ensure reproducibility in our experiments, we use specific random seeds for generating validation and testing instances. For the 10,000 validation instances, we use a random seed of 4321. For the 10,000 testing instances, we use a random seed of 1234. All protocols, including seed selection, align with the practices outlined by [Kool et al. \(2019a\)](#).

B.1.1 TRAVELING SALESMAN PROBLEM (TSP)

The Traveling Salesman Problem (TSP) is a fundamental routing problem that aims to find the Hamiltonian cycle of minimum length. While the original TSP formulation employs mixed-integer linear programming (MILP), in the NCO community, the solution-finding process of TSP is differently formulated for constructive and improvement methods. For constructive methods, the TSP solution is generated by autoregressive solution decoding (i.e., the construction process) in line with [Kool et al. \(2019a\)](#). In each step of node selection, we preclude the selection of nodes already picked in previous rounds. This procedure ensures the feasibility of constructed solutions and also allows for the potential construction of an optimal solution for any TSP instance. For improvement methods, it starts with an initial solution and iteratively searches for an optimal one using local search. In each step, the solution is locally adjusted based on a specified local search operator. We support two representative operators for TSP variants, including the 2-opt in line with [Ma et al. \(2021\)](#) and the flexible k-opt in line with [Ma et al. \(2024\)](#). The former selects two nodes in the current solution and reverses the solution segment between them to perform a 2-opt exchange. The latter selects k nodes so that a k-opt is performed. Both methods ensure the feasibility of the solutions by masking invalid actions. The best solution after a set number of iterations is the final output.

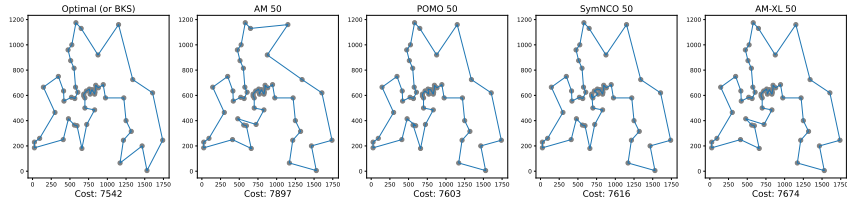


Figure 9: Sample TSP tours on TSPLib’s Berlin 52 with different autoregressive models.

B.1.2 CAPACITATED VEHICLE ROUTING PROBLEM (CVRP)

The Capacitated Vehicle Routing Problem (CVRP) is a popular extension of TSP, applicable to a variety of real-world logistics/routing problems (e.g., delivery services). In CVRP, each node has its own demand, and the vehicle visiting them has a specific capacity and always leaves from a special node called “depot”. The vehicle can visit new nodes while their demand fits in its residual capacity (i.e. the total capacity decreased by the sum of the demands visited in the current path). When no nodes can be added to the path, the vehicle returns to the depot, and its full capacity is restored. Then, it embarks on another tour. The process is repeated until all nodes have been visited. By

applying a similar logic to that of the TSP environment, we can reformulate CVRP as a sequential node selection problem, taking into account demands and capacity.

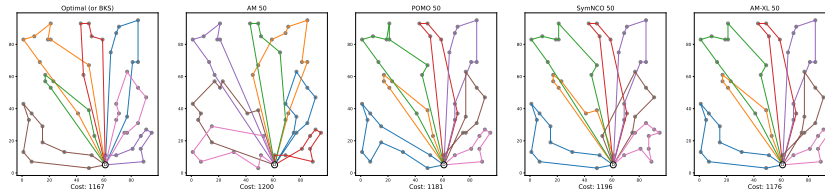


Figure 10: Sample CVRP tours on CVRPLib’s A-n54-k7 instance with different autoregressive models.

Additional generation details To generate the demand, we randomly sample integers between 1 and 10. Without loss of generality, we fix the capacity of the vehicle at 1.0. Then, we normalize the demands by multiplying them by a constant that varies according to the size of the CVRP. The specific constant can be found in our implementation.

B.1.3 ORIENTEERING PROBLEM (OP)

The Orienteering Problem (OP) is a variant of the TSP. In the OP, each node is assigned a prize. The objective of the OP is to find a tour, starting and ending at the depot, that maximizes the total prize collected from visited nodes, while abiding by a maximum tour length constraint. The OP can be framed as a sequential decision-making problem by enforcing the “return to depot” action when no nodes are visitable due to the maximal tour length constraint.

Additional generation details To generate the prize, we use the prize distribution proposed in [Fischetti et al. \(1998\)](#), particularly the distribution that allocates larger prizes to nodes further from the depot.

B.1.4 PRIZE COLLECTING TSP (PCTSP)

In the Prize Collecting TSP (PCTSP), each node is assigned both a prize and a penalty. The objective is to accumulate a minimum total prize while minimizing the combined length of the tour and the penalties for unvisited nodes. By making a minor adjustment to the PCTSP, it can model different subproblems that arise when using the Branch-Price-and-Cut algorithms for solving routing problems.

B.1.5 PICKUP AND DELIVERY PROBLEM (PDP)

The Pickup and Delivery Problem (PDP) is an extension of TSP in the literature [Helsgaun \(2017\)](#); [Ma et al. \(2022\)](#).⁹ In PDP, a pickup node has its own designated delivery node. The delivery node can be visited only when its paired pickup node has already been visited. We call this constraint *precedence constraint*. The objective of the PDP is to find a complete tour with a minimal tour length while starting from the depot node and satisfying the precedence constraints. We assume that *stacking* is allowed, meaning that the traveling agent can visit multiple pickups prior to visiting the paired deliveries. For constructive methods, the PDP solution construction is similar to that of TSP but must obey precedence constraints. For improvement methods, we consider the ruin and repair local search operator presented by [Ma et al. \(2021\)](#). In each step, a pair of pickup and delivery nodes are removed from the current solution and then reinserted back into the solution with potentially better positions. Invalid actions that violate precedence constraints are masked out to ensure the feasibility of PDP solutions.

Additional generation details To generate the positions of the depot, pickups, and deliveries, we sample the node coordinates from the 2D unit square. The first $N/2$ generated nodes are pickups, and the remaining $N/2$ are their respective deliveries. The pickups and deliveries are paired. For a pickup node i , its respective delivery is $i + N/2$ (excluding the depot index).

⁹PDP is also called PDTSP (pickup and delivery TSP).

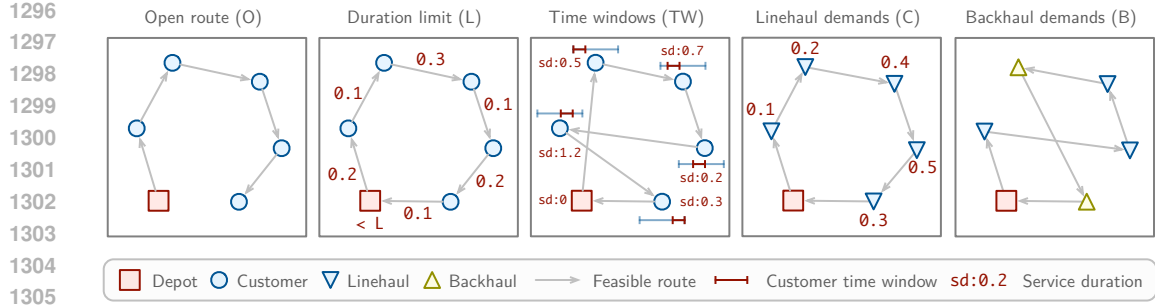


Figure 11: Different VRP attributes. Open routes (O) and duration limits (L) are *global attributes*, whereas time windows (TW), capacitated vehicles for linehaul demands (C) and backhaul demands (B) are *node attributes*. Attributes may be combined in different ways to define VRP variants.

1311 B.1.6 MULTI-TASK VRP (MTVRP)

1312
1313
1314
1315
1316
1317
1318
1319

This environment introduces the 16 VRP variants in Liu et al. (2024a); Zhou et al. (2024) with additional enhancements, such as support for any number of variants in the same batch, as done in Berto et al. (2024). The base logic is the same as CVRP: each node has a demand, and the vehicle has a specific capacity by which it can deliver to nodes and return to the depot to replenish its capacity, with the goal of minimizing the total tour distance. We report each modular constraint definition in the following paragraphs according to Berto et al. (2024); Wouda et al. (2024). Table 8 reports the list of all variants and Fig. 11 illustrates the meaning of each MTRVP component.

1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336

VRP Variant	Capacity (C)	Open Route (O)	Backhaul (B)	Duration Limit (L)	Time Windows (TW)
CVRP	✓				
OVRP	✓	✓			
VRPB	✓		✓		
VRPL	✓			✓	
VRPTW	✓				✓
OVRPTW	✓	✓			✓
OVRPB	✓	✓	✓		
OVRPL	✓	✓		✓	
VRPBL	✓		✓	✓	
VRPBTW	✓		✓		✓
VRPLTW	✓			✓	✓
OVRPBL	✓	✓	✓	✓	
OVRPBTW	✓	✓	✓		✓
OVRPLTW	✓	✓		✓	✓
VRPBLTW	✓		✓	✓	✓
OVRPBLTW	✓	✓	✓	✓	✓

1337
1338
1339
1340

Table 8: The 16 VRP variants that are modeled by the MTRVP environment. All variants include the base Capacity (C). The $k = 4$ features O, B, L, and TW can be combined into any subset, including the empty set and itself (i.e., a *power set*) with $2^k = 16$ possible combinations.

1341
1342
1343

(C) *Demand and Vehicle Capacity* [$q \in [0, Q]$]: Every node i , except the depot, has a demand q_i that must be satisfied by the vehicle with a uniform capacity of $Q > 0$. The sum of the demands served by a vehicle in the same path must not exceed its capacity Q at any point along its route.

1344
1345
1346
1347
1348

(O) *Open Routes* [$o \in \{0, 1\}$]: With open routes, the distance between the last node and the depot is not counted in the total path length. This represents the scenarios where vehicles are not required to return to the depot after serving all assigned customers. Open routes are commonly found in scenarios involving third-party drivers, who are typically compensated only for the deliveries they complete, without the need to return to the depot (Li et al., 2007).

1349

(B) *Backhauls* [$p \in [0, Q]$]: Backhauls extend the concept of demand to include both delivery and pickup requests, thus increasing vehicle utilization and leading to savings. Nodes are categorized as

1350 either linehaul or backhaul nodes.¹⁰ Linehaul nodes require delivery of demand q_i from the depot
 1351 to the node i (similar to CVRP), while backhaul nodes require a pickup of an amount p_i to be trans-
 1352 ported from the node back to the depot. A vehicle can serve both linehaul and backhaul customers in
 1353 a single route, but all linehaul customers must be served before any backhaul customers. A typical
 1354 example of a backhaul problem is a laundry service for hotels that has to deliver clean towels and
 1355 pick up dirty ones, in which the precedence constraint of linehaul nodes is important due to possible
 1356 contamination (Çatay, 2009).

1357 *(L) Duration Limits* [$l \in [0, L]$]: Imposes a limit L on the total travel duration (or distance) of
 1358 each vehicle route, ensuring a fair distribution of workload among different paths. This limit is
 1359 consistently applied to all routes in the problem.

1360 *(TW) Time Windows* [$e, s, l \in [0, T]^3$]: Each node i , except for the depot, has an associated time
 1361 window $[e_i, l_i]$, which specifies the earliest and latest times at which it can be visited. When visiting
 1362 node i , the vehicle must wait for a time s_i before leaving. The vehicle must arrive at customer i
 1363 before the end of its time window l_i , but if they arrive before the start of the time window e_i , they
 1364 must wait at the customer’s location until the time window begins before starting the service. When
 1365 the vehicle returns to the depot, the time is reset to 0.
 1366

1367 **Additional generation details** We introduce the data generation details as follows:

1368 *Locations*: We generate $n + 1$ locations randomly with x_i and $y_i \sim U(0, 1), \forall i \in \{0, \dots, n\}$, where
 1369 $[x_0, y_0]$ represents the depot and $[x_i, y_i], i \in \{1, \dots, n\}$ are the other n nodes.

1370 *Capacity*: The capacity C of the vehicle is determined based on the following calculation:

$$1371 \quad C = \begin{cases} 30 + \lfloor \frac{1000}{5} + \frac{n-1000}{33.3} \rfloor & \text{if } 1000 < n \\ 30 + \lfloor \frac{n}{5} \rfloor & \text{if } 20 < n \leq 1000 . \\ 30 & \text{otherwise} \end{cases}$$

1372
 1373
 1374
 1375 *Open route*: the open route is an instance-wise flag: when set to 1, the route is open, when 0 is
 1376 closed. We sample the flag from a uniform distribution with the same probability of the route being
 1377 open or closed.
 1378

1379 *Linehaul and Backhaul demands*: We generate demands according to the following schema:

- 1380 1. Generate linehaul demands $q_i \in \{0, \dots, Q\}$ for all nodes $i \in \{1, \dots, n\}$. These are needed
 1381 for both backhaul and linehaul scenarios.
- 1382 2. Generate backhaul demands $p_i \in \{0, \dots, Q\}$ for all nodes $i \in \{1, \dots, n\}$.
- 1383 3. For each node $i \in \{1, \dots, n\}$, there is a probability of 0.2 that it is assigned a backhaul
 1384 demand, otherwise, its backhaul demand is set to be 0.
 1385
 1386

1387 Note that even in a backhaul setting, usually not all nodes are backhaul nodes, i.e., we need to
 1388 consider both linehaul and backhaul demands in backhaul problem settings. All demands, both
 1389 linehauls and backhauls, are scaled to $[0, 1]$ through division by the vehicle capacity.

1390 *Duration limits*: Each route is assigned a fixed duration limit L with a default value of 3. We check
 1391 that $2 * d_{0i} < L$ to make sure there is a feasible route for any customer.

1392 *Time Windows*: We generate the time windows for each node $i \in \{1, \dots, n\}$ according to the
 1393 following steps:
 1394

- 1395 1. Generate service time $s_i \in [0.15, 0.18]$.
- 1396 2. Generate time window length $t_i \in [0.18, 0.2]$.
- 1397 3. Calculate distance d_{0i} from node to depot.
- 1398 4. Calculate the upper bound for the start time $h_i = \frac{t_{max} - s_i - t_i}{d_{0i}} - 1$, where t_{max} is the
 1399 maximum time with a default value of 4.6.
 1400

1401 ¹⁰Note that another name of this problem, as adopted in LKH3 (Helsgaun, 2017), is VRP with Pickup and
 1402 Deliveries (VRPPD). However, we align with PyVRP (Wouda et al., 2024) and do not use this name to prevent
 1403 confusion with the *one-to-one PDP*, as we described before, where there is strict precedence between each pair
 of pickup and delivery.

5. Calculate the start time as $e_i = (1 + (h_i - 1) \cdot u_i) \cdot d_{0i}$ with $u_i \sim U(0, 1)$.
6. Calculate the end time as $l_i = e_i + t_i$.

Classical solvers We employ the SotA HGS implementation in PyVRP (Wouda et al., 2024) and OR-Tools (Perron and Furnon, 2023). We make these solvers conveniently available through the `solve` API of the environment.

B.2 SCHEDULING

Scheduling problems are a fundamental class of problems in operations research and industrial engineering, where the objective is to optimize the allocation of resources over time. These problems are critical in various industries, such as manufacturing, computer science, and project management. Currently, RL4CO implements three central scheduling problems, namely the flexible flow shop (FFSP), the job shop (JSSP), and the flexible job shop problem (FJSSP). Each of these problems has unique characteristics and complexities that need to be translated into the environment classes that we will describe hereafter.

B.2.1 JOB SHOP SCHEDULING PROBLEM (JSSP)

The job shop scheduling problem is a well-known combinatorial optimization problem. It is widely used in the operations research community as well as many industries, such as manufacturing and transportation. In the JSSP, a set of jobs J must be processed by a set of machines M . Each job $J_i \in J$ consists of a set of n_i operations $O_i = \{o_{ij}\}_{j=1}^{n_i}$ which must be processed one after another in a given order. The goal of the JSSP is to construct a valid schedule that adheres to the precedence order of the operations and minimizes the makespan, i.e., the time until the last job is finished. One example of such a schedule is shown in Fig. 12.

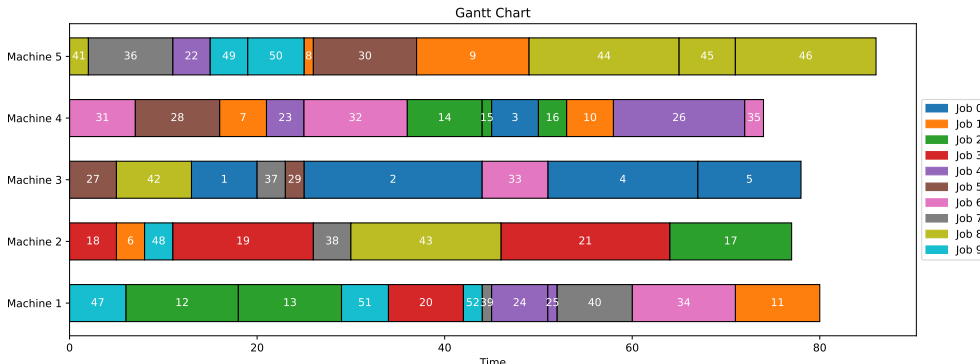


Figure 12: Example Schedule for the JSSP

We formulate the JSSP as a sequential decision problem following the implementation of Tassel et al. (2021). Here, the environment iterates through distinct time steps $t = 1, \dots, T$. At each time step, the agent decides for each machine whether and which job to process next until all machines are busy or all jobs are being processed. In this case, the environment transitions to the next time step at which a machine becomes idle.

Instance Generation We follow the instance generation method described by Zhang et al. (2020), which assumes that each job has exactly one operation per machine, i.e. $n_i = |M|$. Further, processing times for all operations are sampled iid. from a uniform distribution, with parameters specified in Table 9.

B.2.2 FLEXIBLE JOB SHOP SCHEDULING PROBLEM (FJSSP)

The flexible job shop scheduling problem is very similar to the JSSP. However, while in the classical JSSP, each operation $o_{ij} \in O$ has a specified machine and processing time p_{ij} , the flexible job shop

scheduling problem (FJSSP) relaxes this assumption by allowing each operation to be processed by multiple eligible machines $M_k \subseteq M$, potentially with different processing times p_{ijk} associated with the respective operation-machine pair. As a consequence, the agent does not only need to decide which job to process next, but also on which machine it should be processed.

Instance Generation We follow the instance generation method described by Song et al. (2022), who sample n_i operations for each job J_i from a uniform distribution. Further, an average processing time \bar{p}_{ij} is drawn for each operation $o_{ij} \in O$, and the actual processing time per eligible operation-machine pair is subsequently sampled from $U(0.8 \cdot \bar{p}_{ij}, 1.2 \cdot \bar{p}_{ij})$. The parameters used for instance generation can be found in Table 9.

Table 9: Instance generation parameters

	JSSP				FJSSP			
	6×6	10×10	15×15	20×20	10×5	20×5	15×10	20×10
$ J $	6	10	15	20	10	20	15	20
$ M $	6	10	15	20	5	5	10	10
n_i	6	10	15	20	U(4, 6)	U(4, 6)	U(8, 12)	U(8, 12)
\bar{p}_{ij}	U(1, 99)	U(1, 99)	U(1, 99)	U(1, 99)	U(1, 20)	U(1, 20)	U(1, 20)	U(1, 20)
$ M_i $	1	1	1	1	U(1, 5)	U(1, 5)	U(1, 10)	U(1, 10)

B.2.3 FLEXIBLE FLOW SHOP PROBLEM (FFSP)

The flexible flow shop problem (FFSP) is a complex and widely studied optimization problem in production scheduling. It involves N jobs to be processed in S stages, each containing multiple machines ($M > 1$). Each job must pass through the stages in a specified order, but within each stage, it can be processed by any available machine. A critical constraint is that no machine can process more than one job at a time. The objective is to find an optimal schedule that minimizes the total time required to complete all jobs. We formulate the FFSP as a sequential decision process, where at each time step $t = 0, 1, \dots$ and for each idle machine, the agent must decide whether and which job to schedule. If all machines are busy or all jobs are currently being processed, the environment moves to the next time step $t + 1$, and the process repeats until all jobs for each stage have been scheduled.

Instance Generation We follow the data generation process described by Kwon et al. (2021), who sample processing times for each job-machine pair and for every stage independently from a discrete uniform distribution.

B.3 ELECTRONIC DESIGN AUTOMATION

c. This involves solving complex problems that can be either continuous, such as cell placement (Hou et al., 2024), or combinatorial, like decap placement (Kim et al., 2023). RL4CO integrates CO problems in EDA as benchmarking environments.

B.3.1 DECAP PLACEMENT PROBLEM (DPP)

The decap placement problem (DPP) is an electronic design automation problem (EDA) in which the goal is to maximize the performance with a limited number of the decoupling capacitor (decap) placements on a hardware board characterized by asymmetric properties, measured via a probing port. The decaps cannot be placed on the location of the probing port or in keep-out regions (which represent other hardware components) as shown in Fig. 13. The optimal placement of a given number of decaps can significantly impact electrical performance, specifically in terms of power integrity (PI) optimization. PI optimization is crucial in modern chip design, including AI processors, especially with the preference for 3D stacking memory systems like high bandwidth memory (HBM) (Hwang et al., 2021). For comprehensive details, we follow the configuration guidelines provided in (Kim et al., 2023).

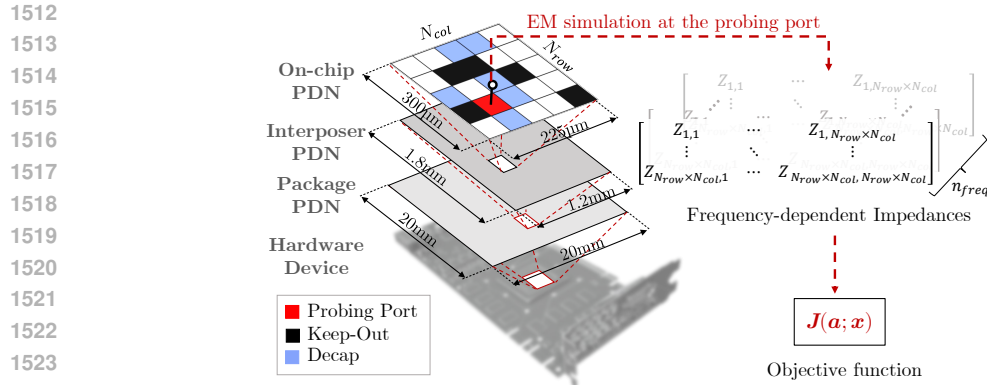


Figure 13: Grid representation of the target on-chip PDN for the DPP problem with a single probing port from Kim et al. (2023).

Baseline solvers We employ two meta-heuristic baselines commonly used in hardware design as outlined in (Kim et al., 2023): random search (RS) and genetic algorithm (GA) (Juang et al., 2021). GA has shown promise as a method for addressing the decap placement problem (DPP).

Instance generation details We use the same data for simulating the hardware board as Kim et al. (2023), with power distribution network (PDN) datasets from Park et al. (2023a). We randomly select one probing port and a number between 1 and 50 keep-out regions sampled from a uniform distribution for generating instances. As in the routing benchmarks, we select seed 1234 for testing the 100 instances.

B.3.2 MULTI-PORT DECAP PLACEMENT PROBLEM (MDPP)

We further consider a more complex and realistic version compared to Kim et al. (2023). The multi-port decap placement problem (mDPP) is a generalization of DPP from Appendix B.3.1 in which measurements from multiple probing ports are performed. The objective function can be either the mean of the reward from the probing ports: 1) (*Maxsum*): the objective is to maximize the average PI among multiple probing ports and 2) (*Maxmin*): maximize the minimum PI between them.

Instance generation details The generation details are the same as DPP, except for the probing port. A number of probing ports between 2 and 5 is sampled from a uniform distribution, and probing ports are randomly placed on the board, just like the other components.

B.4 GRAPH

Many CO problems can be (re-)formulated on graphs (Khalil et al., 2017). In typical CO problems on graphs, actions are defined on nodes/edges, while problem variables and constraints are incorporated in graph topology and node/edge attributes (e.g., weights). The graph-based formulation gives us concise and systematic representations of CO problems. Moreover, existing traditional and machine-learning algorithms for graphs are off-the-shelf tools.

B.4.1 FACILITY LOCATION PROBLEM (FLP)

The optimal usage of limited resources is an important problem to consider in many different fields and has various forms. One specific form of such a problem can be formulated as the facility location problem (FLP), where one aims to choose a given number of locations among given candidates, and the objective is to minimize the overall cost of service (e.g., the sum of the distance from the users to the nearest facility) (Drezner and Hamacher, 2004).

Many real-world problems can be abstracted as instances of FLP. For example, franchise brands may need to determine where to open new retail stores to maximize accessibility and profitability (Shan et al., 2019); governments may need to consider the placement of public facilities (e.g., hospitals

and schools) to maximize the convenience for citizens to use them (Marianov et al., 2002); energy companies may need to determine the best locations for power centers (e.g., power plants and wind farms) to minimize transmission losses (Lotfi et al., 2018).

Formal definition We consider the following specific form of the facility location problem (FLP) used in existing NCO literature (Wang et al., 2022; Bu et al., 2024): (1) given a group of n locations $x_1, x_2, \dots, x_n \in \mathbb{R}^d$ in a d -dimensional space (usually $d = 2$ or 3) and $k < n$, (2) we aim to choose k locations $x_{i_1}, x_{i_2}, \dots, x_{i_k}$ among the given n locations as the locations of facilities, (3) to minimize the sum of the distance from all the n locations to the nearest facility, i.e., $\sum_{j=1}^n \min_{t=1}^k \text{dist}(x_j, x_{i_t})$. We specially consider the Euclidean distance, i.e., $\text{dist}(x_i, x_j) = \|x_i - x_j\|_2$.

Instance generation details The locations are ($d = 2$)-dimensional generated i.i.d. at random. For each location, each coordinate is sampled i.i.d. uniformly at random between 0 and 1. Each instance contains $n = 100$ locations, and $k = 10$ locations are to be chosen.

Classical solvers We apply two MIP solvers, Gurobi (Gurobi Optimization, 2021) and SCIP (Bestuzheva et al., 2021), to obtain the optimal solutions.

B.4.2 MAXIMUM COVERAGE PROBLEM (MCP)

In many real-world scenarios, one needs to allocate limited resources to achieve maximum coverage, which is a fundamental concern across various domains. One specific formulation is called the maximum coverage problem (MCP), where the goal is to select a subset of sets from a given family of sets to maximize the coverage, i.e., the (weighted) size of the union of the selected sets (Khuller et al., 1999).

As a mathematical abstraction, the MCP can be used to represent many real-world problems. For example, radio frequency identification (RFID) system engineers may need to set RFID readers in an optimal way to ensure the maximum coverage of RFID tags (Ali et al., 2011); marketers may need to choose proper forms of advertisement to reach the maximum number of customers (Sun et al., 2018); in security applications (e.g., deploying security cameras), one may need to select the optimal deployment to maximize the coverage of the protected area (Murray et al., 2007).

Formal definition We consider the following specific form of the maximum coverage problem (MCP) used in existing NCO literature (Wang et al., 2022; Bu et al., 2024): (1) given m items (WLOG, $[m] := \{1, 2, 3, \dots, m\}$), where each item t has weight w_t , and a family of n sets $S_1, S_2, \dots, S_n \subseteq [m]$ for some positive integer m and $k < n$, (2) we aim to choose k sets $S_{i_1}, S_{i_2}, \dots, S_{i_k}$ among the given n sets, (3) to maximize the total weighted coverage of the k chosen sets, which is the sum of the weights of items contained in any chosen set, i.e., $\sum_{t \in \cup_{j=1}^k S_{i_j}} w_t$.

Instance generation details First, $m = 200$ items are generated, and the item weights are generated i.i.d., where each weight is a random integer sampled between 1 and 10 (inclusive) uniformly at random. Then, $n = 100$ sets are generated i.i.d., where for each set, we first sample its size between 5 and 15 uniformly at random and then choose that number of items uniformly at random. After generation, $k = 10$ locations are to be chosen.

Classical solvers We apply two MIP solvers, Gurobi (Gurobi Optimization, 2021) and SCIP (Bestuzheva et al., 2021), to obtain the optimal solutions.

B.5 ADDITIONAL ENVIRONMENTS AND BEYOND

We also include in the library additional environments that have been implemented but not fully benchmarked in this paper yet, such as the ATSP, mTSP, Skill-VRP, SMTWTP, and SPCTSP, to name a few. We did not count these in the total environment count (hence the “conservative” estimate). Moreover, several projects, among which co-authors of this paper, have adapted several new environments to their own tasks, which may be included in the future.

Although RL4CO already contains several environments, we acknowledge that the library can be further extended within new directions, which we briefly describe. One such direction is multi-objective combinatorial optimization (Lin et al., 2022; Chen et al., 2024), which is a recently trending research topic of practical importance. Moreover, providing modular reward evaluators to optimize different objectives (for instance, min-max, tardiness) is another avenue of research that we recommend exploring (Park et al., 2023b). Of practical importance is also non-euclidean routing, which so far has received comparatively less attention in this field but is practically important (i.e., DIMACS challenge¹¹). Finally, multi-agent CO (Falkner and Schmidt-Thieme, 2020; Tang et al., 2024a;b; Bettini et al., 2023) is another interesting area of research, which recent approaches model as a sequential decision-making process (Son et al., 2024; Zheng et al., 2024).

Implementing new environments is relatively easy: we created a notebook under the `examples/` folder demonstrating how one can implement a custom environment from the base logic to a fully functioning model. We expect to host an even wider variety of environments in the future, thanks to the community, and invite contributors to help us in our journey.

C BASELINES

This section provides an overview of the key components and methods implemented in RL4CO that can be used as baselines for comparative evaluation. The term “baselines” broadly refers to both the RL algorithms that define the learning objectives and update rules, as well as the policy architectures that parameterize the agent’s behavior in the environment, given that several papers introduce a mix of RL training schemes and policy improvements. We categorize baselines into:

1. **General-purpose RL algorithms (C.1)**
2. **Constructive autoregressive (AR) methods (C.2)**
3. **Constructive non-autoregressive (NAR) methods (C.3)**
4. **Improvement methods (C.4)**
5. **Active search methods (C.5)**

C.1 GENERAL-PURPOSE RL ALGORITHMS

In the following descriptions of RL algorithms, we use the notations of a full problem instance \mathbf{x} and a complete solution \mathbf{a} for simplicity. However, note that these algorithms are also applicable to the usual notion of the sum of rewards over partial states s_t and actions a_t .

C.1.1 REINFORCE (SUTTON ET AL., 1999)

REINFORCE (also known as policy gradients in the literature) is an online RL algorithm whose loss function gradient is given by:

$$\nabla_{\theta} \mathcal{L}_a(\theta|\mathbf{x}) = \mathbb{E}_{\pi(\mathbf{a}|\mathbf{x})} [(R(\mathbf{a}, \mathbf{x}) - b(\mathbf{x})) \nabla_{\theta} \log \pi(\mathbf{a}|\mathbf{x})], \quad (5)$$

where $b(\cdot)$ is a baseline function used to stabilize training and reduce gradient variance. The choice of $b(\cdot)$ can greatly influence the final performance.

C.1.2 ADVANTAGE ACTOR-CRITIC (A2C) (KONDA AND TSITSIKLIS, 1999)

A2C is an algorithm that can be used to solve the RL objective in Eq. (3). It consists of an actor (policy network) and a critic (value function estimator). The actor is trained to maximize the expected cumulative reward by following the policy gradient, while the critic is trained to estimate the value function. The advantage function, computed as the difference between the reward $R(\mathbf{a}, \mathbf{x})$ and the value function $V(\mathbf{x})$, is used to weight the policy gradient update for the actor. This can be seen as a modification of the REINFORCE gradient, where the baseline $b(\mathbf{x})$ is replaced by the value function $V(\mathbf{x})$:

$$\nabla_{\theta} \mathcal{L}_a(\theta|\mathbf{x}) = \mathbb{E}_{\pi(\mathbf{a}|\mathbf{x})} [(R(\mathbf{a}, \mathbf{x}) - V(\mathbf{x})) \nabla_{\theta} \log \pi(\mathbf{a}|\mathbf{x})]. \quad (6)$$

¹¹<http://dimacs.rutgers.edu/programs/challenge/vrp/>

The critic is updated by minimizing the mean-squared error between the estimated value function and the target value, which is the reward for the given problem instance \mathbf{x} :

$$\mathcal{L}_c = \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} (R(\mathbf{a}, \mathbf{x}) - V(\mathbf{x}))^2. \quad (7)$$

By using the advantage function, A2C reduces the variance of the policy gradient and stabilizes training compared to the standard REINFORCE algorithm.

C.1.3 PROXIMAL POLICY OPTIMIZATION (PPO) (SCHULMAN ET AL., 2017)

PPO is another algorithm that can be used to solve the RL objective in Eq. (3). It is an on-policy algorithm that aims to improve the stability of policy gradient methods by limiting the magnitude of policy updates. To this end, PPO introduces a surrogate objective function that constrains the probability ratio between the target policy π_θ that is optimized and a reference policy $\pi_{\theta_{\text{old}}}$, which is periodically updated. This clipping mechanism prevents drastic changes to the target policy, ensuring more reliable and stable learning. Formally, the PPO objective function is given by:

$$\mathcal{L}_{\text{CLIP}}(\theta) = \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} \left[\mathbb{E}_{\mathbf{a} \sim \pi_{\theta_{\text{old}}}(\mathbf{a}|\mathbf{x})} \left[\min \left(\frac{\pi_\theta(\mathbf{a}|\mathbf{x})}{\pi_{\theta_{\text{old}}}(\mathbf{a}|\mathbf{x})} A^{\pi_{\theta_{\text{old}}}(\mathbf{x}, \mathbf{a})}, \right. \right. \right. \\ \left. \left. \left. \text{clip} \left(\frac{\pi_\theta(\mathbf{a}|\mathbf{x})}{\pi_{\theta_{\text{old}}}(\mathbf{a}|\mathbf{x})}, 1 - \epsilon, 1 + \epsilon \right) A^{\pi_{\theta_{\text{old}}}(\mathbf{x}, \mathbf{a})} \right) \right] \right], \quad (8)$$

where θ_{old} represents the parameters of the reference policy, typically a periodically created copy of the parameters θ of the target policy. Further, $A^{\pi_{\theta_{\text{old}}}(\mathbf{x}, \mathbf{a})}$ is the advantage function estimated using the reference policy, and ϵ is a hyperparameter that controls the clipping range, typically set to a small value like 0.2.

The advantage function in PPO is estimated using a learned value function $V_\phi(\mathbf{x})$, where ϕ represents the parameters of the value function. The advantage is computed as:

$$A^{\pi_{\theta_{\text{old}}}(\mathbf{x}, \mathbf{a})} = R(\mathbf{a}, \mathbf{x}) - V_\phi(\mathbf{x}). \quad (9)$$

The value function is learned by minimizing the mean-squared error between the estimated value and the actual return:

$$\mathcal{L}_V(\phi) = \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} \left[(R(\mathbf{a}, \mathbf{x}) - V_\phi(\mathbf{x}))^2 \right]. \quad (10)$$

An optimization step in PPO updates both, the parameters θ of the target policy and the parameters ϕ of the value function by combining $\mathcal{L}_{\text{CLIP}}$ and $\mathcal{L}_V(\phi)$ in a single loss $\mathcal{L}_{\text{PPO}} = \mathcal{L}_{\text{CLIP}} + \beta \mathcal{L}_V(\phi)$, where β is a hyperparameter Schulman et al. (2017).

C.2 CONSTRUCTIVE AUTOREGRESSIVE (AR)

C.2.1 ATTENTION MODEL (AM) (KOOL ET AL., 2019A)

The Attention Model (AM) from Kool et al. (2019a) is an encoder-decoder architecture based on the self-attention mechanism Vaswani et al. (2017) that is at the heart of several state-of-the-art NCO methods, including RL-based ones (Kwon et al., 2020; Kim et al., 2022; Hottung et al., 2024) as well as (self-)supervised ones (Drakulic et al., 2023; Luo et al., 2024a;b). In the original AM, only node features are considered: with abuse of notation from Fig. 3, we consider the `InitEmbedding` as the *node embedding*, and split the *context embedding* into a `ContextEmbedding` which updates the current query and `DynamicEmbedding` that updates the current cached keys and values.

Multi-Head Attention Before delving into the encoder and decoder structures, we briefly introduce the notion of Multi-Head Attention (MHA) from Vaswani et al. (2017), since it is used across several NCO methods. MHA allows the model to jointly attend to information from different representation subspaces at different positions, enabling it to capture various relationships between the input elements. Importantly, it is flexible in handling a variable number of elements.

In the MHA operation, the input sequences Q (queries), K (keys), and V (values) are linearly projected to H different subspaces using learned matrices W_i^Q , W_i^K , and W_i^V , respectively, where H is the number of attention heads:

$$Q_i = QW_i^Q \quad (11)$$

$$K_i = KW_i^K \quad (12)$$

$$V_i = VW_i^V \quad (13)$$

for $i = 1, \dots, H$.

The attention weights are computed as the scaled dot-product between the queries and keys, followed by a softmax operation:

$$A_i = \text{Softmax} \left(\frac{Q_i K_i^T}{\sqrt{d_k}} + M \right) \quad (14)$$

where d_k is the dimension of the keys, used as a scaling factor to prevent the dot-products from getting too large, and M is an optional mask matrix that can be used to prevent attention to certain positions (e.g. infeasible actions in a CO problem).

The output of each attention head is computed as the weighted sum of the values, using the attention weights:

$$Z_i = A_i V_i \quad (15)$$

Finally, the outputs of all attention heads are concatenated and linearly projected using a learned matrix W^O to obtain the final output of the MHA operation:

$$\text{MHA}(Q, K, V) = \text{Concat}(Z_1, \dots, Z_H) W^O \quad (16)$$

This multi-head attention mechanism allows the model to learn different attention patterns and capture various dependencies between the input elements, enhancing the representational power of the model. The queries, keys, and values can come from the same input sequence (self-attention, i.e. $Q = K = V$) or from different sequences (cross-attention), depending on the application. While the attention operation is at the core of much of the current SotA deep learning (Touvron et al., 2023), this scales as $O(L)^2$ where L is the sequence length, such as the number of nodes in a TSP. Thus, an efficient implementation such as FlashAttention (Dao et al., 2022; Dao, 2023) is important, as shown in Appendix E.7.2.

Encoder The encoder’s primary task is to encode input \mathbf{x} into a hidden embedding \mathbf{h} . The structure of f_θ comprises two trainable modules: the `InitEmbedding` and encoder blocks. The `InitEmbedding` module typically transforms problem features into the latent space and problem-specific compared to the encoder blocks, which often involve plain multi-head attention (MHA):

$$\mathbf{h} = f_\theta(\mathbf{x}) \triangleq \text{EncoderBlocks}(\text{InitEmbedding}(\mathbf{x})) \quad (17)$$

Each encoder block in the AM is composed of an Attention Layer, similar to Vaswani et al. (2017). Each layer ℓ is composed of multi-head attention (MHA) for message passing and a Multi-Layer Perceptron (MLP, also known as *feed-forward network (FFN)*), with skip-connections and normalization (Norm):

$$\hat{\mathbf{h}} = \text{Norm} \left(\mathbf{h}^{(\ell-1)} + \text{MHA}(\mathbf{h}^{(\ell-1)}, \mathbf{h}^{(\ell-1)}, \mathbf{h}^{(\ell-1)}) \right) \quad (18)$$

$$\mathbf{h}^{(\ell)} = \text{Norm} \left(\hat{\mathbf{h}} + \text{MLP}(\hat{\mathbf{h}}) \right) \quad (19)$$

with $\ell = [1, \dots, N]$ where N is the number of encoding layers and $\mathbf{h}^0 = \text{InitEmbedding}(\mathbf{x})$. In the encoder side, we have $Q = K = V = \mathbf{h}^{(\ell-1)}$, hence self-attention.

The original implementation of the AM uses $N = 3$ layers $H = 8$ heads of dimension $d_k = \frac{d_h}{M} = 16$, an MLP with one hidden layer of dimension 512 with a ReLU activation function, and a Batch Normalization (Ioffe and Szegedy, 2015) as normalization.

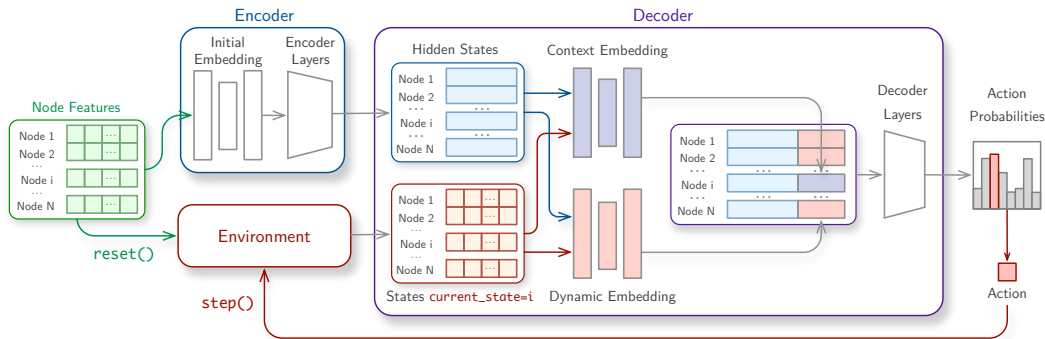


Figure 14: An overview of the modularized Attention Model policy in RL4CO.

Decoder The decoder g_θ autoregressively constructs the solution based on the encoder output \mathbf{h} and the state at current step t , s_t . The solution decoding involves iterative steps until a complete solution is constructed: at each step, starting from the current node’s i query q_t^i

$$q_t^i = \text{ContextEmbedding}(\mathbf{h}, s_t), \quad (20)$$

$$h_t^c = \text{MHA}(q_t^i, K_t^g, V_t^g, M_t), \quad (21)$$

$$\mathbf{z} = \frac{V_t^p h_t^c}{\sqrt{d_k}} \quad (22)$$

where M_t is the set of feasible actions (i.e. the `action_mask`), projections $K_t^g, V_t^g, V_t^p = W_k^g \mathbf{h}, W_v^g \mathbf{h}, W_v^p \mathbf{h}$ can either be precomputed once as cache or updated via a dynamic embedding $K_t^g, V_t^g, V_t^p = \text{DynamicEmbedding}(W_k^g \mathbf{h}, W_v^g \mathbf{h}, W_v^p \mathbf{h}, s_t, \mathbf{h}, \mathbf{x})$, depending on the problem. We note that Eq. (22) is usually referred to as the pointer mechanism (in the codebase, we refer to Eq. (21) and Eq. (22) as the `PointerAttention`). Finally, logits \mathbf{z} (unnormalized output of policy π) are transformed into a probability distribution over the action space:

$$p = \text{Softmax}(C \cdot \tanh(\mathbf{z})) \quad (23)$$

where logits \mathbf{z} for infeasible actions can be set to $-\infty$ to avoid choosing them; and the C value (called *tanh clipping*, usually set to 10) serves in improving the exploration (Bello et al., 2017). We note that Eq. (23) can also include additional operations such as temperature scaling, top-k, and top-p filtering.

Baseline Kool et al. (2019a) additionally introduces the *rollout* baseline b for Eq. (5). At the end of each epoch, a greedy rollout of a baseline policy π_{BL} is executed for each of the sampled instances \mathbf{x} , whose values become baselines for REINFORCE. The algorithm compares the current training policy with a saved baseline policy (similar to the DQN target network (Mnih et al., 2015)) at the end of every epoch, and replace the parameters of π_{BL} with the current trained π if the improvement is significant with a paired t-test of (i.e., 5% in the original paper).

C.2.2 PTR-NET (VINYALS ET AL., 2015)

The original Pointer Network (Ptr-Net) is introduced in Vinyals et al. (2015) and further refined to be trained with RL in (Bello et al., 2017). The base architecture predates the AM (Kool et al., 2019a): an attention mechanism is employed to select outputs of variable length, thus “pointing” at them. The baseline architecture additionally uses an LSTM (Hochreiter and Schmidhuber, 1997), which in practice has less expressivity than full-fledged attention.

C.2.3 POMO (KWON ET AL., 2020)

POMO introduces the *shared* baseline to lower the REINFORCE variance. The key idea is that one can sample rollouts when decoding by forcing diverse starting nodes, which is a powerful inductive bias for certain problems, such as the TSP, in which multiple optimal initial starting points exist. The baseline b_{shared} is the average of all rollouts:

$$b_{\text{shared}}(s) = \frac{1}{N} \sum_{j=1}^N R(\mathbf{a}_j, \mathbf{x}) \quad (24)$$

where N is the number of sampled trajectories (typically set as the number of nodes).

C.2.4 SYMNCO (KIM ET AL., 2022)

SymNCO considers the symmetric nature of combinatorial problems and solutions. There are two major symmetries in combinatorial optimization: 1) *Problem symmetries*: The representation of the input 2D coordinates should have equivalent optimal solution sets and 2) *Solution symmetries*: Multiple permutations can represent an identical cyclic line graph. To reflect this symmetric nature, SymNCO augments the AM architecture by incorporating an auxiliary invariant representation loss function to ensure input 2D symmetries. Additionally, SymNCO employs a shared baseline as Eq. (24) similar to POMO but samples rollouts from both different symmetric problem inputs and solutions together. The implementation is not vastly different from AM and POMO; the primary addition is the symmetric-aware augmentation functions.

1836 C.2.5 POLYNET (HOTTUNG ET AL., 2024)

1837
1838 The PolyNet method proposed by Hottung et al. (2024) enables the learning of a set of complemen-
1839 tary solution strategies within a single model. This facilitates the easy sampling of diverse solutions
1840 at test time, resulting in improved exploration of the search space and, consequently, enhanced over-
1841 all performance. Unlike many other approaches, PolyNet does not artificially increase exploration
1842 by forcing diverse starting actions, as initially proposed by Kwon et al. (2020). Instead, PolyNet uti-
1843 lizes its inherent diversity mechanism, based on its novel architecture and the Poppy loss (Grinsztajn
1844 et al., 2023; Chalumeau et al., 2024):

$$1845 \nabla_{\theta} \mathcal{L} = \mathbb{E}_{\pi(\mathbf{a}^*|\mathbf{x})} [(R(\mathbf{a}^*, \mathbf{x}) - b_o(\mathbf{x})) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}^*|\mathbf{x})], \quad (25)$$

1846 to facilitate exploration during the search process, where \mathbf{a}^* is the *best* solution of K PolyNet sam-
1847 ples and $b_o(\mathbf{x})$ is the average reward of the K samples. This can improve performance for problems
1848 in which the first action greatly influences the performance.

1849 C.2.6 HAM (LI ET AL., 2021)

1851 The Heterogeneous Attention Model (HAM) (Li et al., 2021) is a model specialized for Pickup
1852 and Delivery problems (PDP, Appendix B.1.5), characterized by hard one-to-one precedence con-
1853 straints. To differentiate between pickup and delivery pairs, it introduces *ad hoc* encoder blocks with
1854 a specialized attention mechanism that can differentiate between pickup and delivery pairs.

1856 C.2.7 MTPOMO (LIU ET AL., 2024A)

1857 The MTPOMO developed by Liu et al. (2024a) proposes to adopt a unified model to learn across
1858 various VRP variants. It is motivated by the fact that the diverse VRPs are different combinations of
1859 several shared underlying attributes. By training on a limited number of VRPs with basic attributes,
1860 the model is capable of generalizing to a vast array of VRP variants, each representing different
1861 combinations of these attributes. This approach extends POMO (Kwon et al., 2020) by incorporating
1862 an attribute composition block, facilitating learning across different problems. The cross-problem
1863 learning demonstrates promising zero-shot generation performance on unseen VRPs and benefits
1864 out-of-distribution performance.

1866 C.2.8 MVMoE (ZHOU ET AL., 2024)

1867 The MVMoE architecture proposed by Zhou et al. (2024) incorporates mixture-of-experts
1868 (MoEs) (Jacobs et al., 1991; Jordan and Jacobs, 1994; Shazeer et al., 2017) into attention-based
1869 model (e.g., POMO (Kwon et al., 2020)), such that the model capacity can be greatly enhanced with-
1870 out a proportional increase in computation. For the *encoder* part, MVMoE replaces a feed-forward
1871 network (FFN) with an MoE layer, which typically consists of 1) m experts $\{E_1, E_2, \dots, E_m\}$,
1872 each of which is also an FFN with independent trainable parameters, and 2) a gating network G
1873 parameterized by W_G , which decides how the inputs are distributed to experts. Given a single input
1874 x , $G(x)$ and $E_j(x)$ denote the output of the gating network (i.e., an m -dimensional vector), and the
1875 output of the j th expert, respectively. The output of an MoE layer is calculated as:

$$1876 \text{MoE}(x) = \sum_{j=1}^m G(x)_j E_j(x). \quad (26)$$

1877
1878 The gating algorithm follows the node-level input-choice gating proposed by Shazeer et al. (2017),
1879 which leverages a sparse gating network: $G(x) = \text{Softmax}(\text{Top}K(x \cdot W_G))$. In this way, only
1880 k experts with partial model parameters are activated, hence saving the computation. For the *de-*
1881 *coder* part, MVMoE replaces the final linear layer of MHA with an MoE layer, including m linear
1882 layers and a gating network G . To balance the empirical performance and computational com-
1883 plexity, a hierarchical gating mechanism is further proposed to utilize MoEs during decoding effi-
1884 ciently. In this case, the MoE layer in the decoder includes two gating networks $\{G, G'\}$, m experts
1885 $\{E_1, E_2, \dots, E_m\}$, and a dense layer D . Given a batch of inputs X , the hierarchical gating routes
1886 them in two stages. In the first stage, G' decides to distribute inputs X to either the sparse or dense
1887 layer. In the second stage, if X is routed to the sparse layer, the gating network G is activated to
1888 route nodes to experts on the node level by using the default gating algorithms, i.e., the input-choice
1889 gating. Otherwise, X is routed to the dense layer D and transformed into $D(X)$. In summary, the
hierarchical gating learns to output $G'(X)_0 \sum_{j=1}^m G(X)_j E_j(X)$ or $G'(X)_1 D(X)$. Empirically, hi-

erarchical gating has been found to be more efficient, albeit with a slight sacrifice in in-distribution performance, while demonstrating superiority with out-of-distribution data.

C.2.9 L2D (ZHANG ET AL., 2020)

Learning to Dispatch (L2D) proposed by Zhang et al. (2020) is a DRL method to solve the JSSP. It comprises of the usual encoder-decoder structure, where a graph convolution network (GCN) is employed to extract hidden representations from the JSSP instance. To this end, L2D formulates the JSSP as a disjunctive graph, with nodes reflecting the operations of the problem instance. Nodes of operations that belong to the same job are connected via directed arcs, specifying their precedence relation. Moreover, operations to be processed on the same machine are connected using undirected arcs. Using the resulting neighborhood \mathcal{N} of the nodes, the GCN performs message passing between adjacent operations to construct their hidden representations. Formally, let \mathbf{h}^0 be the initial embeddings of operations O and $\tilde{\mathbf{A}}$ the adjacency matrix with added self-loops of operations, then a graph convolutional layer can be described as follows:

$$\mathbf{h}^{(l+1)} = \sigma \left(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{h}^{(l)} \mathbf{W}^{(l)} \right)$$

Here, $\mathbf{h}^{(l)}$ are the operation embeddings at layer l , $\mathbf{W}^{(l)}$ is a trainable weight matrix at layer l , and $\sigma(\cdot)$ is an activation function such as ReLU. Further, $\tilde{\mathbf{D}}$ is the diagonal degree matrix of $\tilde{\mathbf{A}}$, ensuring appropriate scaling of the features.

Given the operation embeddings, the decoder of L2D first extracts for each job the embedding of the operation that needs to be scheduled next and then feeds them to an MLP $f : \mathbb{R}^{J \times d} \rightarrow \mathbb{R}^{J \times 1}$ to obtain logits for each job $j \in (1, \dots, J)$. In contrast to Kool et al. (2019a) for example, who encode the CO problem once and then generate actions autoregressively using only the decoder, Zhang et al. (2020) use the GCN encoder after each step to generate new hidden representations that reflect the current state of the problem.

C.2.10 HGNN (SONG ET AL., 2022)

The heterogeneous graph neural network (HGNN) is a neural network architecture proposed by Song et al. (2022) to solve the FJSSP. Similar to L2D, HGNN considers an FJSSP instance as a graph. However, instead of treating an FJSSP instance as a disjunctive graph, Song et al. (2022) formulate it as heterogeneous graph with operations and machines posing different node types. Again, operations are connected to each other via directed arcs that specify the precedence relation. Machines are only connected to operations that they are able to process, and the edge weights indicate the respective processing times. To encode the graph, HGNN first projects operations $O \in x$ and machines $M \in x$ into a mutual embedding space \mathbb{R}^d using type-specific transformations \mathbf{W}^O and \mathbf{W}^M , respectively. Given the initial hidden representations \mathbf{h}_i^0 and \mathbf{h}_k^0 for operations $o_i \in O$ and machines $m_k \in M$, respectively, as well as edge embeddings \mathbf{h}_{ik} , an HGNN layer conducts weighted message passing between operations and machines using the processing times of operation-machine pairs:

$$\mathbf{h}_i^{l+1} = \sum_{j \in \mathcal{N}_i} \epsilon_j \mathbf{h}_j^l, \quad \text{where} \quad (27)$$

$$\epsilon_{ij} = \text{Softmax}(\mathbf{a}^\top [\mathbf{h}_j^l || \mathbf{h}_{ij}]). \quad (28)$$

Since operations in the FJSSP can be processed by multiple machines, the decoder must specify not only which job to process next but also on which machine the operation of the selected job should be executed. To this end, Song et al. (2022) concatenates the hidden representations of every operation with the embeddings of every machine. The resulting embeddings are fed to an MLP $f : \mathbb{R}^{J \times M \times 2d} \rightarrow \mathbb{R}^{J \times M \times 1}$, which generates the sampling probabilities for the respective action.

C.2.11 MATNET (KWON ET AL., 2021)

The MatNet architecture proposed by Kwon et al. (2021) adjusts the attention model Kool et al. (2019a) so that it is applicable to bipartite graphs with node types \mathcal{I} and \mathcal{J} as well as a weight matrix $E \in \mathbb{R}^{|\mathcal{I}| \times |\mathcal{J}|}$ corresponding to the edges connecting nodes from the two sets. The novelty of this architecture is that instead of using self-attention as in the attention model, MatNet uses cross-

attention to perform message passing between both node sets and augments the resulting attention scores with the weight matrix E . Formally, let \mathcal{Z} be the set of all nodes $i \in \mathcal{I} \cup \mathcal{J}$, \mathcal{Z}_{ϕ_i} the subset of nodes of the same type as i and $\mathcal{Z}_{\phi_i}^{\mathcal{C}}$ the set of nodes of the respective type. Then, cross-attention is defined as:¹²

$$\alpha'_{ij} = \frac{\mathbf{q}_i^\top \mathbf{k}_j}{\sqrt{d_k}}, \quad \forall i \in \mathcal{Z}, j \in \mathcal{Z}_{\phi_i}^{\mathcal{C}} \quad (29)$$

where

$$\mathbf{q}_i = W_{\phi_i}^Q \mathbf{h}_i^{l-1} \quad \mathbf{k}_j = W_{\phi_i}^K \mathbf{h}_j^{l-1} \quad (30)$$

and weight matrices $W_{\phi_i}^Q$ and $W_{\phi_i}^K \in \mathbb{R}^{d_k \times d_h}$ being learned by the update function corresponding to nodes of type ϕ_i . After that, MatNet augments α'_{ij} with the corresponding edge weight e_{ij} and maps it through a feed-forward neural network $\text{FF} : \mathbb{R}^2 \rightarrow \mathbb{R}$ to a scalar score, which is then normalized using the softmax function:

$$\alpha_{ij} = \frac{\exp(\epsilon_{ij})}{\sum_{q \in \mathcal{Z}_{\phi_i}^{\mathcal{C}}} \exp(\epsilon_{iq})}, \quad \epsilon_{ij} = \text{FF}([\alpha'_{ij} || e_{ij}]) \quad (31)$$

The resulting weights are used to compute a weighted average of the embeddings $\mathbf{v}_j = W_{\phi_i}^V \mathbf{h}_j^{l-1}$ of the nodes in $\mathcal{Z}_{\phi_i}^{\mathcal{C}}$. In the end, skip connections, layer normalization (LN), and feed-forward layers are used as in Vaswani et al. (2017). Besides the original MatNet implementation, RL4CO also implements a version that applies both self- and cross-attention, successively as proposed by Luttmann and Xie (2024). This makes MatNet not only applicable to bipartite graph problems but to the more general class of heterogeneous graphs Luttmann and Xie (2024).

C.2.12 DEVFORMER (KIM ET AL., 2023)

We employ online RL variants of DevFormer (Kim et al., 2023) (DF), an Attention-Model (Kool et al., 2019a) variant specifically designed for autoregressive construction of DPP solutions from Appendix B.3.1. We note that the DF training scheme was initially designed for offline training; however, in this study, we benchmark DF as a sample-efficient online reinforcement learning approach. We benchmark the DF version for RL with the same node and context embedding structure as the original in Kim et al. (2023). We modify the embeddings in the mDPP environment (Appendix B.3.2) version to include the location of multiple probing ports. Min-max and min-sum mDPP versions utilize the same embeddings and are trained separately.

C.3 CONSTRUCTIVE NON-AUTOREGRESSIVE (NAR)

C.3.1 DEEPACO (YE ET AL., 2023)

Ant Colony Optimization (ACO) is an evolutionary algorithm that has been successfully applied to various COPs. Traditionally, customizing ACO for a specific problem requires the expert design of knowledge-driven heuristics. However, this routine of algorithm customization exhibits certain deficiencies: 1) it requires extra effort and makes ACO less flexible; 2) the effectiveness of the heuristic measure heavily relies on expert knowledge and manual tuning; and 3) designing a heuristic measure for less-studied problems can be particularly challenging, given the paucity of available expert knowledge.

DeepACO is designed to automatically strengthen the heuristic measures of existing ACO algorithms and dispense with laborious manual design in future ACO applications. DeepACO consists of two stages: 1) training a neural model to map a COP instance to its heuristic measures, and 2) incorporating the learned heuristic measures into ACO to bias solution constructions and local search. During the training phase, DeepACO parameterizes the heuristic space with a graph neural network (GNN) (Joshi et al., 2019). It trains the GNN across COP instances with REINFORCE, towards minimizing the expected objective value of both constructed solutions and solutions refined by local search. During the inference phase, DeepACO utilizes the well-trained GNN to generate heuristic

¹²For succinctness, note that we omit head and layer enumeration.

measures for ACO. Optionally, DeepACO interleaves local search with neural-guided perturbation to refine the constructed solutions. For more details, please refer to (Ye et al., 2023).

DeepACO is the first NAR model implemented in RL4CO, laying the foundation for other NAR models later integrated into RL4CO. DeepACO offers a versatile methodological framework that allows for further algorithmic enhancements in neural architecture, training paradigms, decoding strategies, and problem-specific adaptations. Notable improvements over DeepACO are introduced by GFACS (Kim et al., 2024).

C.3.2 GFACS (KIM ET AL., 2024)

While DeepACO (Ye et al., 2023) provides promising results and opens new doors for pretraining heuristic measures for the ACO algorithm using deep learning, their method is sub-optimal for two major reasons. Firstly, they utilized policy gradient reinforcement learning (RL), which is an on-policy method that cannot leverage powerful off-policy techniques such as local search. Secondly, their method cannot effectively capture the multi-modality of heuristic distribution because the RL method cannot accurately model multi-modal probabilistic distributions considering the symmetric nature of combinatorial space, where multiple trajectories can lead to identical solutions.

The methodology of GFACS shares a very similar structure with DeepACO. The key difference lies in the learning procedure; GFACS employs generative flow networks (GFlowNets) (Bengio et al., 2021a; 2023) for learning the heuristic matrix. Additionally, they leverage effective off-policy exploration methods using local search. The inference procedure with the learned heuristic matrix remains exactly the same. With the RL4CO modular implementation, both DeepACO and GFACS can run similarly and be comparable at the modular level, allowing future researchers to improve certain modules of training or inference.

C.3.3 GLOP (YE ET AL., 2024B)

Most NCO methods struggle with real-time scaling-up performance; they are unable to solve routing problems involving thousands or tens of thousands of nodes in seconds, falling short of the needs of modern industries. GLOP (Global and Local Optimization Policies) is proposed to address this challenge. It partitions a large routing problem into sub-TSPs and further partitions potentially large (sub-)TSPs into small Shortest Hamiltonian Path Problems (SHPPs). It is the first hybrid method to integrate NAR policies for coarse-grained problem partitions and AR policies for fine-grained route constructions, leveraging the scalability of the former and the meticulousness of the latter.

1) AR (Sub-)TSP Solver. The (Sub-)TSP Solver in GLOP initializes TSP tours using a Random Insertion heuristic, which greedily inserts nodes to minimize cost. These tours are then improved through a process of decomposition and reconstruction. Specifically, the solver decomposes a complete tour into several subtours, which are treated as instances of the Shortest Hamiltonian Path Problem (SHPP). Each subtour is solved using an AR local policy referred to as a “reviser”. These revisers are applied in rounds called “revisions” to enhance the initial tour iteratively. The subtours are normalized and optionally rotated to improve the model’s performance. After solving the SHPP instances, the subtours are reassembled into an improved complete tour. This method allows for efficient and parallelizable improvements on large-scale TSPs.

2) NAR General Routing Solver. The general routing solver in GLOP additionally implements an NAR global policy that either partitions all nodes into multiple sub-TSPs (e.g., for CVRP) or subsets all nodes to form a sub-TSP (e.g., for PCTSP). The NAR global policy is parameterized by a graph neural network (GNN) that processes sparsified input graphs and outputs a partition heatmap. GLOP clusters or subsets nodes by sequentially sampling nodes based on the partition heatmap while adhering to problem-specific constraints. The sub-TSPs are then solved by the (Sub-)TSP solver. The global policy is trained using REINFORCE to output partitions that could lead to the best-performing final solutions after solving sub-TSPs.

GLOP is integrated into RL4CO as the first hybrid method that combines NAR and AR policies, indicating the versatility of RL4CO in accommodating various methodological paradigms. It is promising to further investigate the emerging possibilities that arise when viewing AR and NAR

2052 methods from a unified perspective and combining them synergistically. RL4CO provides a flexible
 2053 and extensible platform for exploring such hybridization in future research.
 2054

2055 C.4 IMPROVEMENT METHODS 2056

2057 Improvement methods leverage RL to train a policy that iteratively performs rewriting exchanges on
 2058 the current solution, aiming to generate a new solution with potentially lower costs. As in construc-
 2059 tive methods, the policy of improvement methods is also based on the encoder-decoder structure.
 2060

2061 C.4.1 DACT (MA ET AL., 2021) 2062

2063 Improvement methods typically take node features and solution features (positional information of
 2064 nodes in the current solution) as key inputs. Encoding VRP solutions involves processing com-
 2065 plex relationships between Node Feature Embeddings (NFEs) and Positional Feature Embeddings
 2066 (PFEs). However, directly adopting the original Transformer to add the two types of embeddings, as
 2067 done by Wu et al. (2021), can cause mixed attention score correlations and impairing performance.
 2068 To address this, the Dual-Aspect Collaborative Transformer (DACT) proposes DAC-Att, which pro-
 2069 cesses NFEs and PFEs separately and employs cross-aspect referential attention to understand the
 2070 consistencies and differences between the two embedding aspects. This approach avoids mixed
 2071 correlations and allows detailed modeling of hidden patterns. Another key issue is the Positional
 2072 Encoding (PE) method. While the original Transformer’s PE works well for linear sequences, it
 2073 may not suit the cyclic nature of VRP solutions. To address this, DACT proposes Cyclic Positional
 2074 Encoding (CPE), inspired by cyclic Gray codes, which generates cyclic real-valued coding vectors to
 2075 capture the topological structure of VRP solutions and improve generalization. Additionally, DACT
 2076 redesigns the RL algorithm for improvement methods, introducing a Proximal Policy Optimization
 with Curriculum Learning (PPO-CL) algorithm to improve training stability and efficiency.

2077 In RL4CO, DACT is implemented and modularized so that other methods can easily reuse com-
 2078 ponents like CPE encoding and the PPO-CL algorithm. It also reuses common parts (such as node
 2079 embedding initialization, decoding functions, etc) from the implementation of constructive methods,
 2080 indicating the flexibility of the RL4CO framework.

2081 C.4.2 N2S (MA ET AL., 2022) 2082

2083 The Neural Neighborhood Search (N2S) method extends the capabilities of improvement methods
 2084 to pickup and delivery problems (PDP). Expanding on the DACT approach, N2S leverages a tai-
 2085 lored MDP formulation for a ruin-repair neighborhood search process. It uses a Node-Pair Removal
 2086 decoder in the ruin stage and a Node-Pair Reinsertion decoder in the repair stage, allowing efficient
 2087 operation on pickup-delivery node pairs. However, more complex decoders increase computational
 2088 costs in the policy network, requiring a balance between encoders and decoders. To address this,
 2089 N2S introduces Synthesis Attention (Synth-Att), which learns a single set of embeddings and synthe-
 2090 sizes attention scores from various node feature embeddings using a Multilayer Perceptron (MLP)
 2091 module. This promotes lightweight policy networks and enhances model expressiveness. The N2S
 2092 encoder with the efficient Synth-Att represents a state-of-the-art design of improvement encoder,
 2093 which is adopted in the latest works Ma et al. (2022; 2024).

2094 In RL4CO, N2S reuses the CPE encoding and the PPO-CL algorithm implemented in DACT. The
 2095 efficient N2S encoder is also modularized and designed to be shared among other improvement
 2096 methods to process the complex relationships between different feature embeddings.
 2097

2098 C.4.3 NEUOPT (MA ET AL., 2024) 2099

2100 A key bottleneck of improvement methods like DACT is their simplistic action space design, which
 2101 typically uses smaller, fixed k values (2-opt or 3-opt) due to decoders struggling with larger, varying
 2102 k . To address this, the latest improvement method introduces Neural k -Opt (NeuOpt), a flexible
 2103 solver capable of handling any given $k \geq 2$. NeuOpt employs an action factorization method to
 2104 break down complex k -opt exchanges into a sequence of basis moves (S-move, I-move, E-move),
 2105 with the number of I-moves determining the k value. This step-by-step construction allows the
 model to automatically determine a suitable k . Similar to variable neighborhood search, NeuOpt
 combines varying k values across search steps, balancing coarse-grained and fine-grained searches,

2106 which is crucial for optimal performance. NeuOpt also features a Recurrent Dual-Stream (RDS)
2107 decoder with recurrent networks and two decoding streams for contextual modeling and attention
2108 computation, effectively capturing the complex dependencies between removed and added edges.

2109 In RL4CO, NeuOpt is implemented by reusing the successful CPE and PPO-CL training modules
2110 from DACT, as well as the efficient encoder from N2S. This demonstrates the strength and versatility
2111 of the RL4CO coding library, which allows for the easy integration of proven methodologies.
2112

2113 C.5 ACTIVE SEARCH METHODS

2114 Active search methods are examples of *transductive* RL, in which an RL algorithm is run to finetune
2115 a pre-trained policy on specific test-time instances.
2116

2117 C.5.1 ACTIVE SEARCH (AS) (BELLO ET AL., 2017)

2118 In active search proposed by Bello et al. (2017), a model is fine-tuned to a single test instance. To
2119 this end, active search uses the same loss formulation as during the original training of the model.
2120 Over the course of the search process, the model’s performance on the single test instance improves,
2121 leading to the discovery of high-quality solutions. While active search is easy to implement, as the
2122 search process closely follows the training process, it is often very slow since all model weights are
2123 adjusted for each test instance individually.
2124

2125 C.5.2 EFFICIENT ACTIVE SEARCH (EAS) (HOTTUNG ET AL., 2022)

2126 Efficient active search (EAS), proposed by Hottung et al. (2022), builds upon the idea of active
2127 search and trains a model on a single instance at test time to enable a guided search. However, EAS
2128 only updates a subset of parameters during the search and allows most operations to be performed
2129 in parallel across a batch of different instances. This approach not only reduces computational costs
2130 but also results in a more stable fine-tuning process, leading to an overall improvement in solution
2131 quality.
2132

2133 D BENCHMARKING SETUP

2134 D.1 METRICS

2135 D.1.1 GAP TO BKS

2136 The Gap to Best Known Solution (BKS) is a commonly used metric to evaluate the performance
2137 of optimization algorithms on benchmark instances. It measures the relative difference between the
2138 best solution found by the algorithm and the BKS for a given problem instance. Given a problem
2139 instance i , let \mathbf{a}_i be the objective value of the best solution found by the algorithm, and let \mathbf{a}_i^* be the
2140 objective value of the BKS for that instance. The Gap to BKS for the i -th instance is defined as:
2141

$$2142 \text{Gap to BKS}_i = 100 \times \left(\frac{\mathbf{a}_i - \mathbf{a}_i^*}{\mathbf{a}_i^*} \right) \quad (32)$$

2143 The Gap to BKS is expressed as a percentage, with a value of 0% indicating that the algorithm
2144 has found a solution that matches the BKS. A positive Gap to BKS indicates that the algorithm’s
2145 solution is worse than the BKS, while a negative Gap to BKS (though less common) indicates that
2146 the algorithm has found a new best solution for the instance¹³.
2147

2148 D.1.2 PRIMAL INTEGRAL

2149 The Primal Integral (PI) is a metric that evaluates the anytime performance of optimization algo-
2150 rithms by capturing the trade-off between solution quality and computational time (Berthold, 2013;
2151

2152 ¹³Note that when calculating the gap for a set of instances, one should do an average of gaps, i.e.
2153 $\frac{1}{n} \sum_{i=1}^n \text{Gap to BKS}_i$, instead of calculating the gap of the average $100 \times \frac{\sum \mathbf{a}_i}{\sum \mathbf{a}_i^*}$, which might yield
2154 similar results in some settings but prone to error especially for certain distributions.

Thyssens et al., 2023). It is defined as the area under the curve of the incumbent solution value plotted against time, normalized by the BKS value and the total time budget:

$$PI = 100 \times \left(\frac{\sum_{i=1}^n \mathbf{a}_{i-1} \cdot (t_i - t_{i-1}) + \mathbf{a}_n \cdot (T_{\max} - t_n)}{T_{\max} \cdot \mathbf{a}^*} - 1 \right) \quad (33)$$

where T_{\max} is the total time budget, \mathbf{a}_i is the incumbent solution value at time t_i , and \mathbf{a}^* is the best known solution value. A lower PI percentage indicates better anytime performance. The PI complements other metrics, such as the Gap to BKS, by providing insights into the temporal aspect of an algorithm’s performance, making it particularly useful for assessing anytime algorithms (Jesus et al., 2020).

D.1.3 RUNTIME MEASUREMENT

Runtime normalization Comparing the run-time efficiency of different methods across various hardware configurations can be challenging. In the RL4CO benchmark, we generally run the inference on a single machine; when this is not possible due to resource limitations, we employ the run-time normalization approach based on the *PassMark* hardware rating¹⁴. This approach normalizes time budgets and run times during the evaluation process, allowing for a more equitable comparison of methods. We use the definition of Accorsi et al. (2022); Thyssens et al. (2023) in normalizing: the reference machine combines a single CPU thread and a single GPU, the *PassMark* score s for GPU-based methods is calculated as:

$$s = \frac{1}{2}(\#CPU \cdot CPU_Mark + \#GPU \cdot GPU_Mark) \quad (34)$$

To normalize the solution time from machine 1 to machine 2, we calculate $\tilde{t}_2 = t_1 \frac{s_1}{s_2}$, where t_1 is the solution time on machine 1, s_1 is the *PassMark* score of machine 1, and s_2 is the *PassMark* score of machine 2. Note that in the case of most classical solvers, the GPU_Mark is simply set to 0 due to them running on CPU.

Cross-solver comparisons Another aspect of NCO evaluation that has to be addressed is the fact that evaluation between classical and learned solvers is often done on different devices, namely on (single-threaded) CPUs and GPUs, respectively. Moreover, while multiple instances in NCO can usually be solved in a batch, this is not usually the case for classical solvers. A more correct way is to measure the *per-instance* solution time (which we do on large-scale NAR routing), which is more realistic for real-world applications. For other studies, we employ the standard procedure of NCO of evaluating times on batches as done in the original methods, making sure to compare “apples with apples” (i.e., different NCO approaches are compared with the same settings). We note that while RL4CO focuses on comparisons between NCO solvers and creating an open-source ecosystem for this specific area, future studies (and possibly works in the RL4CO community) may also include comparisons with classical solvers under different conditions, which we recognize as an important research direction.

D.2 HARDWARE & SOFTWARE

D.2.1 HARDWARE

Most experiments (during testing) were carried out on a machine equipped with two AMD EPYC 7542 32-CORE PROCESSOR CPUs with 64 threads each and four NVIDIA RTX A6000 graphic cards with 48 GB of VRAM, of which only one is used during inference. We note that, due to the amount of experiments and contributions, training was performed on a variety of hardware combinations, particularly University clusters. We found RL4CO to be robust and efficient across different combinations of CPU, GPU, and software. Throughout the text, we may report the hardware setting on which testing took place if it differs from the default one. In case different configurations were used or results were reported from previous works, we refer to Appendix D.1.3 for result standardization.

¹⁴*PassMark*: <https://www.passmark.com/> is also used in the 2022 DIMACS challenge: <http://dimacs.rutgers.edu/programs/challenge/vrp/>.

2214 D.2.2 SOFTWARE

2215
 2216 Software-wise, we used Python 3.11 and PyTorch 2.3 (Paszke et al., 2019)¹⁵, most notably due
 2217 to the native implementation of `scaled_dot_product_attention`. Given that most models
 2218 in RL constructive methods for CO generally use attention for encoding states, FlashAttention has
 2219 some boost on the performance (between 5% and 20% saved time depending on the problem size)
 2220 when training is subject to mixed-precision training, which we do for all experiments. During decod-
 2221 ing, the FlashAttention routine is not called since, at the time of writing, it does not support maskings
 2222 other than causal; this could further boost performance compared to older implementations. Refer
 2223 to Appendix A.2 for additional details regarding notable software choices of our library, namely
 2224 TorchRL, PyTorch Lightning, and Hydra.

2225 D.3 HYPERPARAMETERS

2227 D.3.1 COMMON HYPERPARAMETERS

2228
 2229 Common hyperparameters can be found in the `config/` folder from the RL4CO library, which
 2230 can be conveniently loaded by Hydra. We provide yaml-like configuration files below, divided by
 2231 experiments in Listing 1.

2232 D.3.2 CHANGING POLICY COMPONENTS

2233
 2234 We train the models evaluated in Table 3 using the same number of training instances as well as
 2235 identical hyperparameters. Specifically, models are trained for 10 epochs on 2.000 training instances
 2236 using the PPO algorithm with clip range $\epsilon = 0.2$. The training dataset is split into batches of size
 2237 100 to construct the replay buffer. For the PPO optimization we sample mini-batches of size 512
 2238 from the replay buffer until it is empty and repeat this for $\mathcal{R} = 3$ inner epochs. All models use
 2239 an embedding dimension d_h of 256. The number of encoder layers is set to $L = 3$ in each case.
 2240 Further, MatNet and the AM Pointer use $H = 8$ attention heads. The parameters of the models
 2241 are updated using the Adam optimizer with learning rate 10^{-4} . Afterwards, the trained policies
 2242 are evaluated on 1.000 randomly generated test instances. The Hydra config files corresponding
 2243 to this experiment, which also implement the different model architectures, can be found in the
 2244 `config/experiment/scheduling` folder from the RL4CO library

2245 D.3.3 MIND YOUR BASELINE

2246
 2247 We run all models to match the original implementation details under *controlled* settings. In par-
 2248 ticular, we run all models for 250,000 gradient steps with the same Adam (Kingma and Ba, 2014)
 2249 optimizer with a learning rate of 10^{-4} and 0 weight decay. For POMO, we match the original im-
 2250 plementation details of weight decay as 10^{-6} . For POMO, the number of multistarts is the same as
 2251 the number of possible initial locations in the environment (for instance, for TSP50, 50 starts are
 2252 considered). In the case of Sym-NCO, we use 10 as augmentation for the shared baseline; we match
 2253 the number of effective samples of AM-XL to the ones of Sym-NCO to demonstrate the differences
 2254 between models.

2265 ¹⁵During development, we also used beta wheels as well as manually installed version of FlashAttention
 2266 (Dao et al., 2022; Dao, 2023). Note that software version varied in terms of training runs depending on the
 2267 author who ran experiments (e.g. any range of Python and PyTorch as $[3.9, 3.10, 3.11] \times [2.0, 2.1, 2.2, 2.3]$,
 which RL4CO can support out of the box on multiple devices and operating systems.

```

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

```

Example Hydra Configuration

```

1 defaults: # override default configurations under configs/
2   - override /env: tsp.yaml
3   - override /model: am.yaml
4   - override /callbacks: default.yaml
5   - override /trainer: default.yaml
6   - override /logger: wandb.yaml
7
8 # Environment
9 env:
10  generator_params:
11    num_loc: 50
12
13 # RL Algorithm and policy (env passed automatically)
14 model:
15  policy: # override policy parameters to pass to the RL algo
16    _target_: rl4co.models.zoo.am.policy.AttentionModelPolicy
17    embed_dim: 128
18    num_heads: 8
19    num_encoder_layers: 3
20    feedforward_hidden: 128
21    env_name: "${env.name}" # automatically construct env embeddings
22    baseline: "rollout" # REINFORCE baseline
23    batch_size: 512
24    train_data_size: 1_280_000
25    optimizer_kwargs:
26      lr: 1e-4
27
28 # Optional override of checkpoint parameters
29 model_checkpoint:
30   dirpath: ${paths.output_dir}/checkpoints
31   filename: "epoch_{epoch:03d}"
32
33 # Trainer
34 trainer:
35   max_epochs: 100
36   gradient_clip_val: 1.0
37   max_epochs: 100
38   precision: "16-mixed" # allows for FlashAttention
39   strategy: DDPStrategy # efficient for multiple GPUs
40   matmul_precision: "medium" # speeds up calculation
41
42 # Logging
43 logger:
44   wandb:
45     project: "rl4co"
46     name: "am-tsp${env.generator_params.num_loc}"

```

Listing 1: Example `example.yaml` configuration for the AM from the AR routing experiments. Additional parameters are modularized in the actual configs and moved to the other config folders (such as `env/tsp.yaml`) so that a single experiment config is not too cluttered. Running this configuration is simple: placed under `configs/experiments/`, it can be called with `python run.py experiment=example`.

The number of epochs for all models is 100, except for AM-XL (500). We also employ learning rate scheduling, in particular, `MultiStepLR`¹⁶ with $\gamma = 0.1$ on epoch 80 and 95; for AM-XL, this applies on epoch 480 and 495.

PPO for the AM We follow other hyperparameters for REINFORCE baselines. We set the number of mini-epochs to 2, mini-batch size to 512, clip range to 0.2, and entropy coefficient $c_2 = 0.01$.

¹⁶https://pytorch.org/docs/stable/generated/torch.optim.lr_scheduler.MultiStepLR

2322 Interestingly, we found that normalizing the advantage as done in the Stable Baselines PPO2 imple-
 2323 mentation¹⁷ slightly hurt performance, so we set the `normalize advantage` parameter to `False`. We
 2324 suspect this is because the NCO solvers are trained on *multiple* problem instances, unlike the other
 2325 RL applications that aim to learn a policy for a single MDP.
 2326

2327 **Sample Efficiency Experiments** We keep the same hyperparameters as the *mind your baseline*,
 2328 experiments except for the number of epochs and scheduling. We consider 5 independent runs
 2329 that match the number of samples *per step* (i.e., the batch size is exactly the same for all models
 2330 after considering techniques such as the multistart and symmetric baselines). For AM Rollout, we
 2331 employ half the batch size of other models since it requires double the number of evaluations due to
 2332 its baseline.
 2333

2334 **Search Methods Experiments** For these experiments, we employ the same models trained in the
 2335 in-distribution benchmark on 50 nodes. For Active Search (AS), we run 200 iterations for each in-
 2336 stance and an augmentation size of 8. The Adam optimizer is used with a learning rate of 2.6×10^{-4}
 2337 and weight decay of 10^{-6} . For Efficient Active Search, we benchmark EAS-Lay (with an added
 2338 layer during the single-head computation, `PointerAttention` in our code) with the original
 2339 hyperparameters proposed by Hottung et al. (2022). The learning rate is set to 0.0041 and weight
 2340 decay to 10^{-6} . The search is restricted to 200 iterations with dihedral augmentation of 8 as well as
 2341 imitation learning weight $\lambda = 0.013$.

2342 Testing is performed on 100 instances on both TSP and CVRP for $N \in [200, 500, 1000]$, generated
 2343 with the usual random seed for testing 1234.
 2344

2345 D.3.4 GENERALIZATION: CROSS-TASK AND CROSS-DISTRIBUTION

2346 In addition to training on uniformly distributed instances, as is standard for POMO Kwon et al.
 2347 (2020), we further train POMO Kwon et al. (2020) on a mixture of multiple distributions (i.e., the
 2348 exemplar distributions defined in (Bi et al., 2022)) and multiple VRP tasks (i.e., CVRP, OVRP,
 2349 VRPL, VRPB, VRPTW, and OVRPTW, as defined in (Liu et al., 2024a; Zhou et al., 2024; Berto
 2350 et al., 2024)) with fixed problem size $N = 50$, termed as MDPOMO and MTPOMO, respectively.
 2351 Note that all the models in Table 5 undergo training across 10,000 epochs, each with a batch size of
 2352 512 and 10,000 training instances. The other training setups are consistent with the previous work
 2353 (Kwon et al., 2020). The whole training time is within one day. During inference, we evaluate their
 2354 generalization performance on the benchmark datasets in CVRPLib Lima et al. (2014) using greedy
 2355 rollout with $8 \times$ instance augmentation and multiple start nodes following Kwon et al. (2020).
 2356

2357 D.3.5 LARGE-SCALE INSTANCES

2358 The GLOP (Ye et al., 2024b) models’ global policy are trained on random instances of CVRP1K and
 2359 CVRP2K, respectively. Both models are trained for 100 epochs, with each epoch comprising 1000
 2360 instances. To accelerate the training process, random insertion is utilized as the sub-TSP solver.
 2361

2362 For the experiment results presented in Table 6, we evaluate our implementation using the identical
 2363 instances and setup as those utilized in Ye et al. (2024b). The AM revisers involved are directly
 2364 adopted from Ye et al. (2024b). Table 14 reports the generalization performance of the CVRP2K
 2365 model on 100 CVRP10K instances and 24 CVRP20K instances. These test instances are generated
 2366 following the procedure in Nazari et al. (2018), with the capacities fixed to 1000.

2367 D.3.6 COMBINING CONSTRUCTION AND IMPROVEMENT

2368 To test the potential collaboration between constructive and improvement methods (in Appendix E.5
 2369 and Section 5.3), we recorded the performance of improvement methods during inference with initial
 2370 solutions generated either randomly or by leveraging solutions generated greedily by constructive
 2371 methods. This was done for both TSP and PDP with a fixed problem size of $N = 50$. We used
 2372 a test set with 1,000 instances for both TSP and PDP and recorded the runtime for all constructive
 2373 and improvement solvers based on an INTEL XEON GOLD 5317 CPU @ 3.00GHZ and one RTX
 2374 3090 GPU.
 2375

¹⁷<https://stable-baselines.readthedocs.io/en/master/modules/ppo2.html>

For the constructive models to bootstrap improvement, we used the POMO and HAM (i.e. AM with rollout baseline, with HAM (Li et al., 2021) encoder for construction PDP) directly from Appendix D.3.3. Note that these models were trained under controlled settings and could see a further boost in performance with further training. Moreover, while we used simple greedy evaluation, more complex evaluation schemes may be used, such as combining symmetric augmentation, multistart, or advanced sampling techniques as nucleus sampling.

For the improvement models, we used both DACT and NeuOpt (with $K = 4$) for TSP, and the N2S model for PDP. Training for all models was conducted with 200 epochs and 20 batches per epoch, with a batch size of 512 for TSP and 600 for PDP. The n-step and maximum improvement steps for training were set to 4 and 200, respectively. Other hyperparameters such as learning rate, curriculum learning scaler, and gradient norm clip were set as per their original papers.

D.4 DECODING SCHEMES

Due to the limited space in the main paper, we further elaborate on the setup of the decoding schemes (or *strategies*) in this section, shown in Fig. 15.



Figure 15: Inference methods we consider in RL4CO. These can also be combined together, such as greedy multistart with augmentation.

D.4.1 AUGMENTATIONS

In RL4CO, we consider as augmentations any transformation ψ that maps an instance x into an instance x' whose (optimal) solution should be the same or close to the original. Augmentations have been used in various domains, such as computer vision, where, for example, labels are invariant to rotations. Similarly, in Euclidean CO, one can apply the *dihedral transformation* of Table 10 to generate a new instance whose solution is the same as the original one, composed of 4 rotations and 2 flips for a total of $\times 8$ transformation (which is the default used in POMO-based models as Kwon et al. (2020); Liu et al. (2024a); Zhou et al. (2024)). As introduced in Kim et al. (2022), one may additionally use any angle θ to perform a symmetric transformation as follows:

Table 10: Dihedral transformations (Kwon et al., 2020).

$\psi(x, y)$	
(x, y)	(y, x)
$(x, 1-y)$	$(y, 1-x)$
$(1-x, y)$	$(1-y, x)$
$(1-x, 1-y)$	$(1-y, 1-x)$

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \psi(x, y) = \begin{pmatrix} x \cos \theta & -y \sin \theta \\ x \sin \theta & +y \cos \theta \end{pmatrix}$$

where $\theta \in [0, 2\pi]$. Interestingly, we found that, generally, the dihedral augmentation is worse in terms of sample efficiency compared to randomly augmenting by sampling a θ value. We note that other augmentations are possible, including dilation (Bdeir et al., 2022) (i.e., rescaling) and possibly new ones such as *jittering*, which may have a broader application than Euclidean CO.

D.4.2 SAMPLING

In most NCO approaches, sampling is performed by simply increasing the evaluation budget but without additional modifications that can be important for better performance. We include the following techniques in RL4CO: 1) *Sampling with Softmax Temperature*, 2) *Top-k Sampling* and 3) *Top-p Sampling*, visualized in Fig. 16.

Sampling with Softmax Temperature Sampling with softmax temperature is a technique used to control the randomness of the sampling process. The temperature parameter τ is introduced to the softmax function, which converts the logits z into a probability distribution:

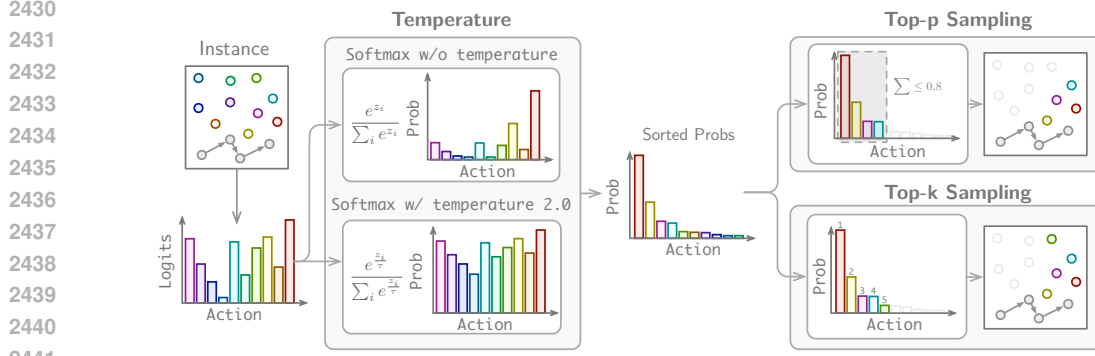


Figure 16: Sampling techniques implemented in RL4CO.

$$p_i = \frac{\exp(z_i/\tau)}{\sum_{j=1}^N \exp(z_j/\tau)} \quad (35)$$

where p_i is the probability of selecting the i -th action, z_i is the corresponding logit, and N is the total number of actions. A higher temperature $\tau > 1$ makes the distribution more uniform, increasing the chances of selecting less likely actions. Conversely, a lower temperature $0 < \tau < 1$ makes the distribution sharper, favoring the most likely actions.

Top-k Sampling Top-k sampling is a method that restricts the sampling space to the k most likely actions. Given the logits z , the top- k actions with the highest probabilities are selected, and the probabilities of the remaining actions are set to zero. The probability distribution is then renormalized over the selected actions:

$$p_i = \begin{cases} \frac{\exp(z_i/\tau)}{\sum_{j \in \mathcal{T}_k} \exp(z_j/\tau)} & \text{if } i \in \mathcal{T}_k \\ 0 & \text{otherwise} \end{cases} \quad (36)$$

where \mathcal{T}_k is the set of indices corresponding to the top- k actions. Top-k sampling helps to eliminate the possibility of generating low-probability actions, improving the quality and coherence of the generated output. We note that, however, in CO problems, it may not be as straightforward as in large language models to select the k parameter since neighborhoods and distributions are not homogeneous.

Top-p Sampling Top-p sampling, also known as nucleus sampling, is an alternative to top-k sampling that dynamically adjusts the number of actions considered for sampling based on a probability threshold p (Holtzman et al., 2019). The actions are sorted by their probabilities in descending order, and the cumulative probability is calculated. The sampling space is then restricted to the smallest set of actions whose cumulative probability exceeds the threshold p :

$$\mathcal{T}_p = \left\{ i : \sum_{j=1}^i p_j \leq p \right\} \quad (37)$$

where \mathcal{T}_p is the set of indices corresponding to the actions included in the top-p sampling. The probabilities of the actions in \mathcal{T}_p are renormalized, while the probabilities of the remaining actions are set to zero:

$$p_i = \begin{cases} \frac{\exp(z_i/\tau)}{\sum_{j \in \mathcal{T}_p} \exp(z_j/\tau)} & \text{if } i \in \mathcal{T}_p \\ 0 & \text{otherwise} \end{cases} \quad (38)$$

Top-p sampling provides a more dynamic way to control the diversity and quality of the generated output compared to top-k sampling. In CO, this is also a more structured way of performing training or evaluation since top-p sampling is agnostic of the number of nodes, unlike top-k sampling.

E ADDITIONAL EXPERIMENTS

E.1 MIND YOUR BASELINE: FURTHER INSIGHTS

Benchmark Setup We focus on benchmarking the AR routing NCO solvers under controlled settings, aiming to compare all benchmarked methods as closely as possible in terms of network architectures and the number of training samples consumed.

Models We evaluate the following NCO solvers: 1) *AM* (Kool et al., 2019a) with rollout baseline, 2) *POMO* (Kwon et al., 2020) with the shared baseline to train AM instead of the rollout baseline; we also use six MHA layers and InstanceNorm instead of BatchNorm according to the original implementation, 3) *Sym-NCO* (Kim et al., 2022) utilizes the symmetric baseline to train AM instead of the rollout baseline and the same encoder as POMO, 4) *AM-XL* is an AM model that adopts *POMO*-style MHA encoder, and trained on the same number of samples as POMO, with the goal of seeing whether training for longer, as done in POMO, can significantly improve the results 5) *A2C*, i.e. AM trained with Advantage Actor-Critic (A2C), 6) *AM-PPO* trained via the Proximal Policy Optimization (PPO, Schulman et al. (2017)) algorithm and finally 7) Polynet (Hottung et al., 2024) with shared baseline and setting $K = n$.

For fairness of comparison, we try to match the number of training steps to be the same and adjust the batch size accordingly. Specifically, we train models for 100 epochs as in Kool et al. (2019a) using the Adam optimizer (Kingma and Ba, 2014) with an initial learning rate (LR) of 0.001 with a decay factor of 0.1 after the 80th and 95th epochs¹⁸. We evaluate the trained solvers using the schemes shown in Fig. 15.

E.1.1 MAIN IN-DISTRIBUTION RESULTS

We first measure the performances of NCO solvers on the same dataset distribution on which they are trained. We first observe that, counter to the commonly known trends that $AM < POMO < Sym-NCO$, the trends can change to decoding schemes and targeting CO problems. Especially when the solver decodes the solutions with *Augmentation* or *Greedy Multistart + Augmentation*, the performance differences among the benchmarked solvers on TSP and CVRP become less significant. Surprisingly, PolyNet performs well even in the greedy one-shot setting, despite its primary focus on generating diverse solutions. For decoding schemes that generate multiple solutions, PolyNet demonstrates strong performance across various problems. Particularly for decoding schemes without multistarts, PolyNet benefits significantly from its inherent diversity mechanism

We note that the original implementation of POMO¹⁹ is not directly applicable to OP, PCTSP, and PDP. Adapting it to solve new problems is not straightforward due to the coupling between environment and policy implementations. However, owing to the flexibility of RL4CO, we successfully implemented POMO for OP and PCTSP. Our results indicate that POMO underperforms in OP and PCTSP; unlike TSP, CVRP, and PDP, where all nodes need to be visited, OP and PCTSP are not constrained to visit all nodes. Due to such differences, POMO’s visiting all nodes strategy may not work as an effective inductive bias. Further, we benchmark the NCO solvers for PDP, which was not originally supported natively by each of the benchmarked solvers. We apply the environment embeddings and the Heterogeneous Attention Encoder from HAM (Li et al., 2021) to the NCO models for encoding pickup and delivery pairs, further emphasizing RL4CO’s flexibility. We observe that AM-XL, which employs the same RL algorithm as AM but features the encoder architecture of POMO and is trained with an equivalent number of samples, yields performance comparable to NCO solvers using more sophisticated baselines. This suggests that careful controls on architecture and the number of training samples are required when evaluating NCO solvers. We also re-implemented

¹⁸We find that simple learning rate scheduling with `MultiStepLinear` can improve performance i.e., compared to the original AM implementation.

¹⁹<https://github.com/yd-kwon/POMO>

Table 11: In-distribution benchmark results for routing problems with 50 nodes. We report the gaps to the best-known solutions of classical heuristics solvers.

Method	TSP			CVRP			OP			PCTSP			PDP		
	Cost ↓	Gap	Time	Cost ↓	Gap	Time	Prize ↑	Gap	Time	Cost ↓	Gap	Time	Cost ↓	Gap	Time
<i>Classical Solvers</i>															
Gurobi	5.70	0.00%	2m	—	—	—	—	—	—	—	—	—	—	—	—
Concorde	5.70	0.00%	2m	—	—	—	—	—	—	—	—	—	—	—	—
HGS	—	—	—	10.37	0.00%	10h	—	—	—	—	—	—	—	—	—
Compass	—	—	—	—	—	—	16.17	0.00%	5m	—	—	—	—	—	—
LKH3	5.70	0.00%	5m	10.38	0.10%	12h	—	—	—	—	—	—	6.86	0.00%	1h30m
OR Tools	5.80	1.83%	5m	—	—	—	—	—	—	4.48	0.00%	5h	7.36	7.29%	2h
<i>Greedy One Shot Evaluation</i>															
A2C	5.83	2.22%	(<1s)	11.16	7.09%	(<1s)	14.77	8.64%	(<1s)	5.15	14.96%	(<1s)	7.52	9.90%	(<1s)
AM	5.78	1.41%	(<1s)	10.95	5.30%	(<1s)	15.46	4.40%	(<1s)	4.59	2.46%	(<1s)	7.51	9.88%	(<1s)
POMO	5.75	0.89%	(<1s)	10.80	3.99%	(<1s)	13.86	14.26%	(<1s)	5.00	11.61%	(<1s)	7.59	10.64%	(<1s)
Sym-NCO	5.72	0.47%	(<1s)	10.87	4.61%	(<1s)	15.67	3.09%	(<1s)	4.52	2.12%	(<1s)	7.39	7.73%	(<1s)
AM-XL	5.73	0.54%	(<1s)	10.84	4.31%	(<1s)	15.69	2.98%	(<1s)	4.53	2.44%	(<1s)	7.31	6.56%	(<1s)
AM-PPO	5.76	0.92%	(<1s)	10.87	4.60%	(<1s)	15.67	3.05%	(<1s)	4.55	2.45%	(<1s)	7.43	8.31%	(<1s)
PolyNet	5.72	0.68%	2s	10.81	4.24%	2s	15.70	2.93%	2s	4.54	2.45%	2s	8.26	3.46%	2s
<i>Sampling with width M = 1280</i>															
A2C	5.74	0.72%	40s	10.70	3.07%	1m24s	15.14	6.37%	48s	4.96	10.71%	57s	7.32	6.70%	1m15s
AM	5.72	0.40%	40s	10.60	2.22%	1m24s	15.90	1.68%	48s	4.52	0.99%	57s	7.25	5.69%	1m15s
POMO	5.71	0.18%	1m	10.54	1.64%	2m30s	14.62	9.56%	1m10s	4.82	7.59%	1m23s	7.31	6.56%	1m50s
Sym-NCO	5.70	0.14%	1m	10.58	2.03%	2m30s	16.02	0.93%	1m10s	4.52	0.82%	1m23s	7.17	4.52%	1m50s
AM-XL	5.71	0.17%	1m	10.57	1.91%	2m30s	15.97	1.25%	1m10s	4.52	0.88%	1m23s	7.15	4.23%	1m50s
AM-PPO	5.70	0.15%	40s	10.52	1.52%	1m24s	16.04	0.78%	48s	4.48	0.18%	57s	7.17	4.52%	1m15s
PolyNet	5.70	0.15%	1m20s	10.42	0.53%	2m40s	16.08	0.52%	1m15s	4.47	0.13%	2m15s	6.93	0.81%	2m10s
<i>Greedy Multistart (N)</i>															
A2C	5.80	1.81%	2s	10.90	4.86%	6s	14.61	9.65%	4s	5.12	14.29%	5s	7.54	9.85%	4s
AM	5.77	1.21%	2s	10.73	3.39%	6s	15.71	2.84%	4s	4.56	1.89%	5s	7.46	8.75%	4s
POMO	5.71	0.29%	3s	10.58	2.04%	8s	13.95	13.71%	7s	4.98	11.16%	7s	7.46	8.75%	6s
Sym-NCO	5.72	0.36%	3s	10.71	3.17%	8s	15.88	1.79%	7s	4.55	1.59%	7s	7.38	7.58%	6s
AM-XL	5.72	0.42%	3s	10.68	2.88%	8s	15.85	1.95%	7s	4.56	1.79%	7s	7.25	5.69%	6s
AM-PPO	5.74	0.61%	2s	10.67	2.72%	6s	15.98	1.21%	4s	4.53	1.18%	5s	7.23	5.39%	4s
PolyNet	5.70	0.25%	3s	10.52	1.42%	18s	16.05	0.71%	3s	4.54	1.31%	10s	7.18	4.65%	5s
<i>Greedy with Augmentation (1280)</i>															
A2C	5.71	0.18%	40s	10.63	2.49%	1m24s	14.89	7.91%	48s	5.15	14.96%	1m	7.03	2.46%	1m15s
AM	5.70	0.07%	40s	10.53	1.56%	1m24s	15.88	1.79%	48s	4.59	2.46%	1m	7.14	4.08%	1m15s
POMO	5.70	0.06%	1m	10.55	1.72%	2m30s	14.23	11.97%	1m15m	5.09	13.61%	1m42s	7.15	4.23%	1m45s
Sym-NCO	5.70	0.01%	1m	10.53	1.54%	2m30s	15.94	1.41%	1m15m	4.58	2.17%	1m42s	7.03	2.48%	1m45s
AM-XL	5.70	0.01%	1m	10.52	1.47%	2m30s	15.90	1.66%	1m15m	4.59	2.54%	1m42s	6.98	1.75%	1m45s
AM-PPO	5.70	0.15%	40s	10.52	1.52%	1m24s	16.01	0.84%	48s	4.48	0.18%	1m	7.00	2.04%	1m15s
PolyNet	5.70	0.17%	1m30s	10.47	0.92%	3m	16.05	0.72%	2m	4.47	0.10%	2m10s	6.94	1.20%	2m15s
<i>Greedy Multistart with Augmentation (N × 16)</i>															
A2C	5.72	0.41%	32s	10.67	2.81%	1m	15.22	5.88%	30s	5.06	12.94%	35s	7.10	3.51%	50s
AM	5.71	0.21%	32s	10.55	1.73%	1m	16.05	0.76%	30s	4.54	1.28%	35s	7.10	3.50%	50s
POMO	5.70	0.05%	48s	10.48	1.11%	2m	15.05	6.94%	1m	4.92	9.81%	1m10s	7.12	3.79%	1m25s
Sym-NCO	5.70	0.03%	48s	10.54	1.63%	2m	16.09	0.51%	1m	4.53	1.17%	1m10s	7.01	2.19%	1m25s
AM-XL	5.70	0.04%	48s	10.53	1.50%	2m	16.08	0.57%	1m	4.54	1.25%	1m10s	7.00	2.04%	1m25s
AM-PPO	5.70	0.03%	32s	10.51	1.45%	1m	16.09	0.49%	30s	4.49	0.89%	35s	6.98	1.75%	50s
PolyNet	5.70	0.15%	1m	10.41	0.36%	2m16s	16.11	0.37%	1m24s	4.49	0.24%	1m35s	7.02	2.33%	1m50s

PointerNetworks (Vinyals et al., 2015; Bello et al., 2017), but we excluded them from the main table due to their poor performance, i.e., more than 4% optimality gap in TSP50.

Table 11 and Table 12 show detailed results for 50 and 20 nodes, respectively.

E.1.2 DECODING SCHEMES COMPARISON

During inference, investing more computational resources (i.e., sampling more), the trained NCO solver can discover improved solutions. We examine the performance gains achieved with varying numbers of samples. As shown in Fig. 17, the Augmentation decoding scheme achieves the Pareto front with limited samples and, notably, generally outperforms other decoding schemes. We note that while sampling with a light decoder can be more efficient in terms of speed than sampling, this may not be true for heavy-decoder (Luo et al., 2024a) or decoder-only models (Drakulic et al., 2023; Luo et al., 2024b; Pirnay and Grimm, 2024), where decoding via greedy augmentations may help improve performance.

E.1.3 SAMPLE EFFICIENCY

We additionally evaluate the NCO solvers based on the number of training samples (i.e., the number of reward evaluations). As shown in Fig. 18, we found that actor-critic methods (e.g., A2C and PPO) can exhibit efficacy in scenarios with limited training samples, as demonstrated by the TSP50/100

Table 12: In-distribution results for models trained on 20 nodes.

Method	TSP			CVRP			OP			PCTSP			PDP		
	Cost ↓	Gap	Time	Cost ↓	Gap	Time	Prize ↑	Gap	Time	Cost ↓	Gap	Time	Cost ↓	Gap	Time
<i>Classical Solvers</i>															
<i>Gurobi</i> [†]	3.84	0.00%	7s	—	—	—	—	—	—	—	—	—	—	—	—
<i>Concorde</i>	3.84	0.00%	1m	—	—	—	5.39	0.00%	16m	3.13	0.00%	2m	—	—	—
<i>HGS</i>	—	—	—	6.13	0.00%	4h	—	—	—	—	—	—	—	—	—
<i>Compass</i>	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
<i>LKH3</i>	3.84	0.00%	15s	6.14	0.16%	5h	—	—	—	—	—	—	—	—	—
<i>OR Tools</i>	3.85	0.37%	1m	—	—	—	—	—	—	3.13	0.00%	5h	4.70	3.16%	1h
<i>CPLEX</i>	—	—	—	—	—	—	—	—	—	—	—	—	4.56	0.00%	7m23s
<i>Greedy One Shot Evaluation</i>															
A2C	3.86	0.64%	(<1s)	6.46	5.00%	(<1s)	5.01	6.70%	(<1s)	3.36	7.35%	(<1s)	4.71	3.31%	(<1s)
AM	3.84	0.19%	(<1s)	6.39	3.92%	(<1s)	5.20	3.17%	(<1s)	3.17	1.28%	(<1s)	4.82	5.70%	(<1s)
POMO	3.84	0.18%	(<1s)	6.33	3.00%	(<1s)	4.69	12.69%	(<1s)	3.41	8.95%	(<1s)	4.85	6.36%	(<1s)
Sym-NCO	3.84	0.05%	(<1s)	6.30	2.58%	(<1s)	5.30	1.37%	(<1s)	3.15	0.64%	(<1s)	4.70	3.07%	(<1s)
AM-XL	3.84	0.07%	(<1s)	6.31	2.81%	(<1s)	5.25	2.23%	(<1s)	3.17	1.26%	(<1s)	4.71	3.29%	(<1s)
PolyNet	3.84	0.10%	(<1s)	6.40	4.44%	(<1s)	5.26	2.28%	(<1s)	3.18	1.98%	(<1s)	4.69	2.92%	(<1s)
<i>Sampling with width M = 1280</i>															
A2C	3.84	0.15%	20s	6.26	2.08%	24s	5.12	4.66%	22s	3.28	4.79%	23s	4.64	1.76%	23s
AM	3.84	0.04%	20s	6.24	1.78%	24s	5.30	1.30%	22s	3.15	0.78%	23s	4.66	2.19%	23s
POMO	3.84	0.02%	36s	6.20	1.06%	40s	4.90	8.83%	37s	3.33	6.39%	39s	4.68	2.63%	39s
Sym-NCO	3.84	0.01%	36s	6.22	1.44%	40s	5.34	0.59%	37s	3.14	0.35%	39s	4.64	1.75%	39s
AM-XL	3.84	0.02%	36s	6.22	1.46%	40s	5.32	0.93%	37s	3.15	0.56%	39s	4.64	1.75%	39s
PolyNet	3.84	0.00%	47s	6.14	0.23%	1m15s	5.35	0.52%	37s	3.13	0.15%	1m15s	4.59	0.57%	1m36s
<i>Greedy Multistart (N)</i>															
A2C	3.85	0.36%	(<1s)	6.33	3.04%	3s	5.06	5.77%	2s	3.30	5.18%	2s	4.85	6.42%	2s
AM	3.84	0.12%	(<1s)	6.28	2.27%	3s	5.24	2.42%	2s	3.16	0.95%	2s	4.67	2.41%	2s
POMO	3.84	0.05%	(<1s)	6.21	1.27%	4s	4.76	11.32%	3s	3.35	7.03%	4s	4.66	2.19%	4s
Sym-NCO	3.84	0.03%	(<1s)	6.22	1.48%	4s	5.32	0.87%	3s	3.15	0.62%	4s	4.69	2.85%	4s
AM-XL	3.84	0.05%	(<1s)	6.22	1.38%	4s	5.29	1.49%	3s	3.15	0.64%	4s	4.65	1.97%	4s
PolyNet	3.84	0.01%	1s	6.17	0.71%	5s	5.34	0.58%	1s	3.15	0.76%	5s	4.81	5.43%	5s
<i>Greedy with Augmentation (1280)</i>															
A2C	3.84	0.01%	20s	6.22	1.35%	24s	5.04	6.10%	22s	3.33	6.39%	23s	4.61	1.11%	23s
AM	3.84	0.00%	20s	6.20	1.07%	24s	5.25	2.25%	22s	3.16	0.96%	23s	4.63	1.54%	23s
POMO	3.84	0.00%	36s	6.18	0.84%	45s	4.85	9.76%	38s	3.37	7.55%	42s	4.62	1.32%	42s
Sym-NCO	3.84	0.00%	36s	6.17	0.71%	45s	5.33	0.77%	38s	3.15	0.63%	42s	4.61	0.95%	42s
AM-XL	3.84	0.00%	36s	6.17	0.68%	45s	5.30	1.30%	38s	3.15	0.68%	42s	4.61	0.96%	42s
PolyNet	3.84	0.00%	55s	6.16	0.48%	1m10s	5.35	0.50%	57s	3.13	0.16%	1m2s	4.59	0.58%	1m10s
<i>Greedy Multistart with Augmentation (N × 16)</i>															
A2C	3.84	0.01%	9s	6.20	1.12%	48s	5.20	3.17%	32s	3.28	4.95%	25s	4.75	4.06%	23s
AM	3.84	0.00%	9s	6.18	0.78%	48s	5.34	0.56%	32s	3.14	0.32%	25s	4.63	1.52%	23s
POMO	3.84	0.00%	13s	6.16	0.50%	1m	5.09	5.29%	45s	3.35	6.95%	38s	4.61	1.10%	42s
Sym-NCO	3.84	0.00%	13s	6.17	0.61%	1m	5.35	0.39%	45s	3.14	0.24%	38s	4.60	0.89%	42s
AM-XL	3.84	0.00%	13s	6.16	0.44%	1m	5.35	0.46%	45s	3.14	0.28%	38s	4.60	0.87%	42s
PolyNet	3.84	0.00%	18s	6.14	0.16%	1m20s	5.37	0.31%	1m	3.13	0.12%	58s	4.61	1.03%	55s

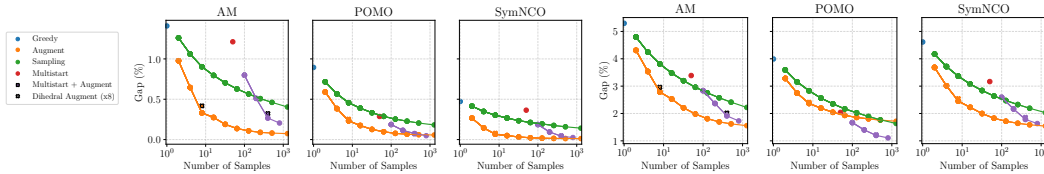


Figure 17: Pareto front of decoding schemes by the number of samples. Left: TSP50; right: CVRP50.

results in Fig. 18. This observation suggests that NCO solvers with control over the number of samples may exhibit a different trend in sample efficiency: if reward function evaluation is expensive, REINFORCE baselines that include additional reward function evaluations such as Greedy Rollout, POMO, and SymNCO may be sample-inefficient. While this is not the case for most CO problems (for instance: in routing, it is inexpensive to calculate routes), in other areas as Electronic Design Automation, where reward evaluation is resource-intensive due to the necessity of electrical simulations, in which sample efficiency can become even more crucial.

E.1.4 OUT-OF-DISTRIBUTION

In this section, we evaluate the out-of-distribution performance of the NCO solvers by measuring the gap compared to the best-known solutions (BKS). The evaluation results are visualized in Fig. 19. Contrary to the in-distribution results, we find that NCO solvers with sophisticated baselines (i.e.,

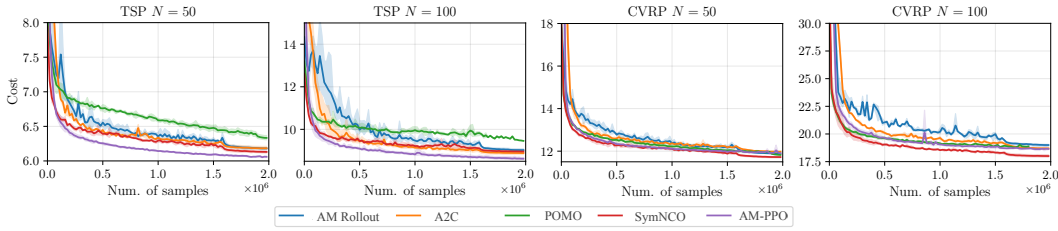


Figure 18: Validation cost curves and number of training samples consumed. Models with greater performance after full training may show worse convergence properties when the number of training samples is limited.

POMO and Sym-NCO) tend to exhibit worse generalization when the problem size changes, either for solving smaller or larger instances. This can be seen as an indication of “overfitting” to the training sizes. On the other hand, variants of AM show relatively better generalization results overall.

Besides, we also evaluate the model by sampling decoding strategy with different temperatures as shown in Fig. 20, k values for Top- k as shown in Fig. 21, and p values for Top- p as shown in Fig. 22. A higher temperature or a lower p value with Top- p sampling can improve the generalization ability on large-scale problems, while Top- k sampling has limited contribution to generalization cross problem sizes.

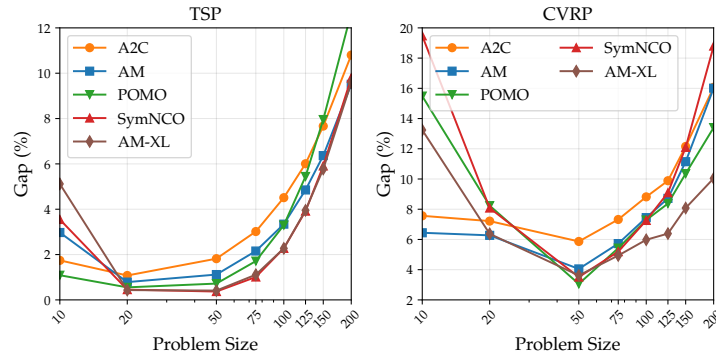


Figure 19: Out-of-distribution generalization by greedy decoding for models with different reinforce baselines trained on 50 nodes. Stronger performance in distribution does not always translate to out-of-distribution.

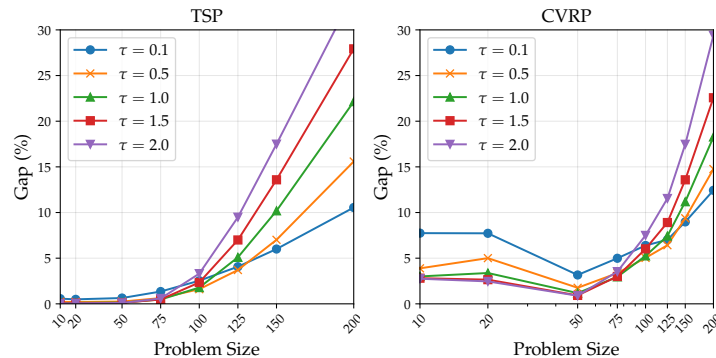


Figure 20: Out-of-distribution generalization by sampling with different temperatures τ for POMO trained on 50 nodes.

E.1.5 SEARCH METHODS

A way to adapt to distribution changes is using *transductive RL*, commonly known as (active) search methods, which involve training (a part of) a pre-trained NCO solver to adapt to CO instances of interest. We evaluate 1) *Active Search (AS)* (Bello et al., 2017) which finetunes a pre-trained model

2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711

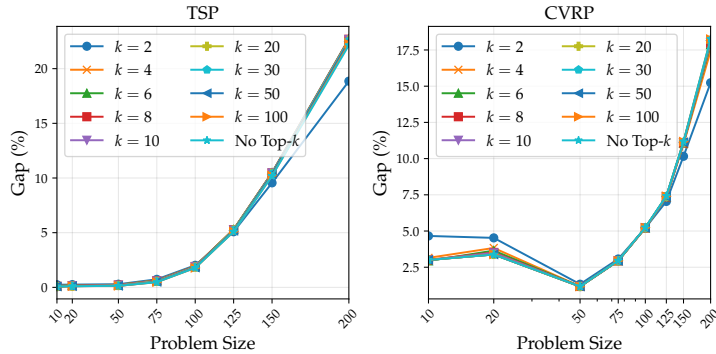


Figure 21: Out-of-distribution generalization by sampling with different Top- k for POMO trained on 50 nodes.

2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726

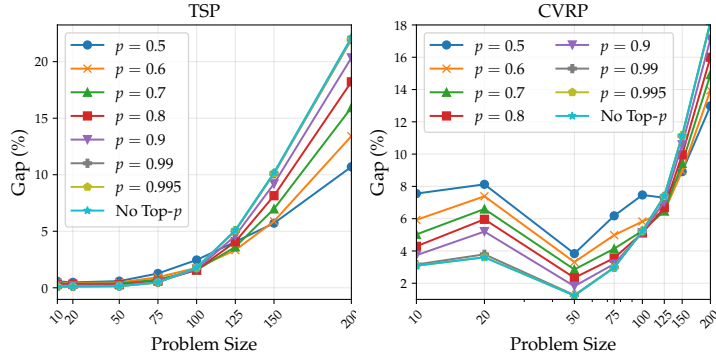


Figure 22: Out-of-distribution generalization by sampling with different Top- p for POMO trained on 50 nodes.

Table 13: Search Methods results of models pre-trained on 50 nodes. *Classic* refers to Concorde (Applegate et al., 2023) for TSP and HGS (Vidal, 2022; Wouda et al., 2024) for CVRP. OOM is "Out of Memory".

Type	Metric	TSP						CVRP					
		POMO			Sym-NCO			POMO			Sym-NCO		
		200	500	1000	200	500	1000	200	500	1000	200	500	1000
<i>Classic</i>	Cost	10.17	16.54	23.13	10.72	16.54	23.13	27.95	63.45	120.47	27.95	63.45	120.47
<i>Zero-shot</i>	Cost	13.15	29.96	58.01	13.30	29.42	56.47	29.16	92.30	141.76	32.75	86.82	190.69
	Gap[%]	29.30	81.14	150.80	24.07	77.87	144.14	4.33	45.47	17.67	17.17	36.83	58.29
	Time[s]	2.52	11.87	96.30	2.70	13.19	104.91	1.94	15.03	250.71	2.93	15.86	150.69
<i>AS</i>	Cost	11.16	20.03	OOM	11.92	22.41	OOM	28.12	63.98	OOM	28.51	66.49	OOM
	Gap[%]	4.13	21.12	OOM	11.21	35.48	OOM	0.60	0.83	OOM	2.00	4.79	OOM
	Time[s]	7504	10070	OOM	7917	10020	OOM	8860	21305	OOM	9679	24087	OOM
<i>EAS</i>	Cost	11.10	20.94	35.36	11.65	22.80	38.77	28.10	64.74	125.54	29.25	70.15	140.97
	Gap[%]	3.55	26.64	52.89	8.68	37.86	67.63	0.52	2.04	4.21	4.66	10.57	17.02
	Time[s]	348	1562	13661	376	1589	14532	432	1972	20650	460	2051	17640

2744
2745
2746
2747
2748
2749
2750
2751
2752
2753

on the searched instances by adapting all the policy parameters and 2) *Efficient Active Search (EAS)*: from (Hottung et al., 2022) which finetunes a subset of parameters (i.e., embeddings or new layers) and adds an imitation learning loss to improve convergence.

We apply AS and EAS to POMO and Sym-NCO pre-trained on TSP and CVRP with 50 nodes to solve larger instances having $N \in [200, 500, 1000]$ nodes. As shown in Table 13, solvers with search methods improve the solution quality. However, POMO generally shows better improvements over Sym-NCO. This suggests once more that the “overfitting” of sophisticated baselines can perform better in training distributions but eventually worse in different downstream tasks.

E.1.6 ADDITIONAL LARGE-SCALE RESULTS

We also show in Table 14 additional large-scale results with $10k+$ nodes obtained with the hybrid AR/NAR GLOP model (Ye et al., 2024b). Fig. 23 demonstrates a solution obtained through our implementation of GLOP for CVRP35K. It represents the maximum scale of CVRP that RL4CO is capable of solving within 24GB of graphics memory while preserving the performance.

Table 14: Performance on large-scale CVRP instances with ten thousands of nodes.

	CVRP10K		CVRP20K	
	Obj.	Time	Obj.	Time
HGS (Vidal, 2022)	108.1	4.01h	182.7	6.03h
Random Insertion	187.9	0.16s	330.4	0.61s
GLOP-G (Insertion)	127.0	2.42s	208.3	10.9s
GLOP-G (AM)	119.6	4.68s	199.6	14.8s
GLOP-G (LKH)	111.4	5.06s	191.4	17.9s

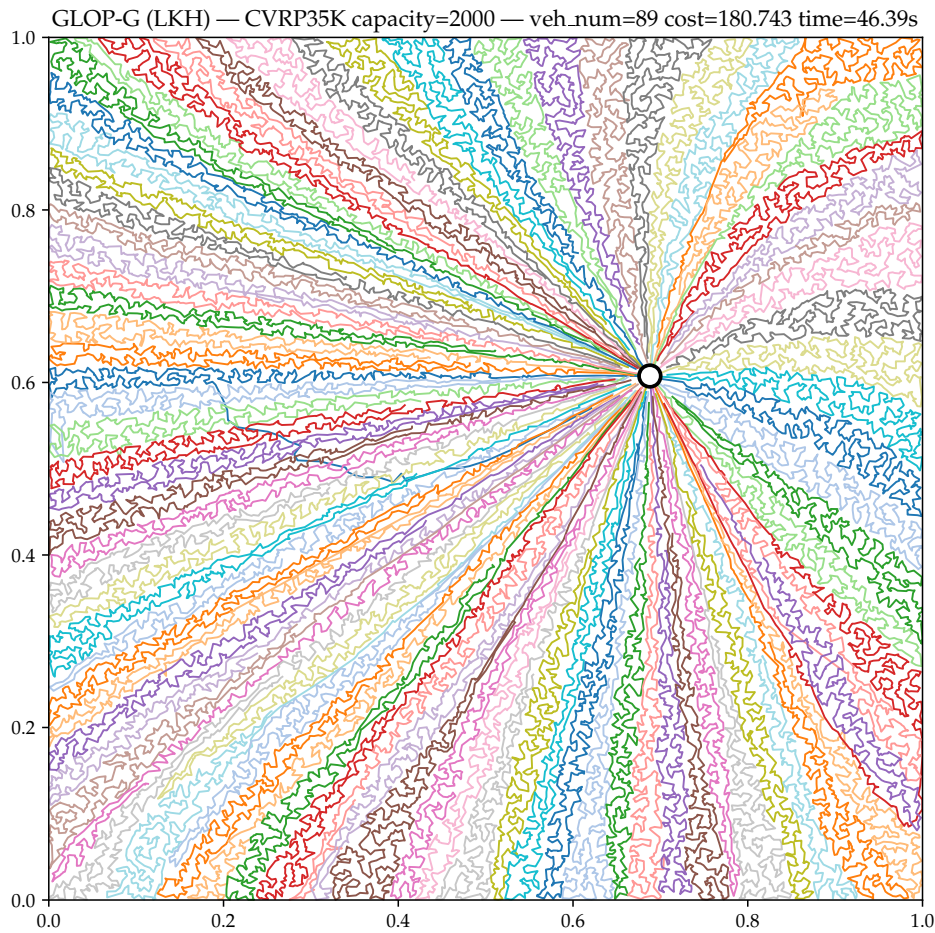


Figure 23: A visualization of the solution generated by GLOP on CVRP35K.

2808 Table 15: Benchmarking results of ACO method in TSP with 200, 500, 1000 nodes. The reported values are
 2809 obtained by averaging over 128 test instances. The time is the average computation time for solving a single
 2810 instance.

Method	TSP200			TSP500			TSP1000		
	Cost	Gap(%)	Time(s)	Cost	Gap(%)	Time(s)	Cost	Gap(%)	Time(s)
<i>Concorde</i> Applegate et al. (2023)	10.72	0.00	0.9	16.55	0.00	10.7	23.12	0.00	108.3
<i>ACO</i>	10.88	1.52	1.0	17.23	4.11	4.0	24.42	5.65	19.8
<i>DeepACO</i>	10.80	0.79	1.0	16.87	1.95	4.3	23.82	3.03	20.7
<i>GFACS</i>	10.75	0.32	1.0	16.80	1.56	4.3	23.78	2.87	20.7

2817 Table 16: Benchmarking results of ACO methods with different τ values in TSP with 500 nodes. The reported
 2818 values are the average cost of 128 test instances.

Method	$\tau = 0.05$	$\tau = 0.1$	$\tau = 0.25$	$\tau = 0.5$	$\tau = 0.75$	$\tau = 1.0$	$\tau = 1.5$	$\tau = 2.0$
<i>ACO</i>	17.05	16.95	17.03	17.11	17.19	17.23	17.26	17.26
<i>DeepACO</i>	17.00	16.97	16.92	16.84	16.85	16.87	16.88	16.89
<i>GFACS</i>	16.92	16.90	16.86	16.80	16.80	16.80	16.81	16.82

2826 E.2 LEARNING HEURISTICS FOR ANT COLONY OPTIMIZATION

2827 E.2.1 EXPERIMENT SETTINGS

2828 We adhered to the hyperparameters specified in the original papers for DeepACO ([Ye et al., 2023](#))
 2829 and GFACS ([Kim et al., 2024](#)) for GFlowNets training. We conducted two distinct benchmarks for
 2830 ACO methods. The first benchmark evaluated the ability to solve the Traveling Salesman Problem
 2831 (TSP) at different scales: 200, 500, and 1000. We use the test instances provided by DeepACO²⁰.
 2832 The second benchmark assessed inference capability at various temperature values of τ in TSP with
 2833 500 nodes. The temperature τ is a hyperparameter for the heatmap distribution of the heuristic
 2834 matrix in ACO, where a low τ emphasizes exploitation and a high τ emphasizes exploration. For
 2835 both experiments, the optimality gaps are calculated with respect to the average cost of solutions
 2836 obtained using Concorde [Applegate et al. \(2023\)](#).
 2837

2838 E.2.2 RESULTS

2839 **TSP Benchmark** [Table 15](#) shows the results for the first benchmark. In this benchmark, we ob-
 2840 served that GFACS outperforms other baselines, and DeepACO surpasses ACO. These results are
 2841 consistent with their respective claims ([Ye et al., 2023](#); [Kim et al., 2024](#)), providing evidence that
 2842 our benchmark is sufficiently valid. Notably, our algorithm also performed slightly faster than the
 2843 original implementation, likely due to the batchified environment of RL4CO.
 2844

2845 **Performance Comparison for Different Heatmap Temperatures (τ)** [Table 16](#) shows the re-
 2846 sults for the second benchmark. This benchmark compared inference performance across different
 2847 heatmap temperatures (τ). We observed notable performance variation with changes in τ . This
 2848 highlights the importance of inference and sampling strategies even after deep network training is
 2849 completed. Additionally, GFACS produced more consistent results with different τ values. This
 2850 provides empirical evidence of the robustness of GFACS, which is due to its ability to model a sam-
 2851 pler capable of generating diverse and high-reward solutions. The modularization of RL4CO allows
 2852 for a focused study on inference capabilities, enabling future researchers to contribute to this aspect
 2853 using the RL4CO pipeline.
 2854

2855 E.3 LEARNING TO SCHEDULE

2856 Compared to routing problems, scheduling problems have not been extensively studied by the NCO
 2857 community. On the one hand side, NCO methods for scheduling are harder to benchmark due to
 2858 the absence of well-performing heuristics like the LKH algorithm for the TSP. On the other hand,
 2859 scheduling problems involve more complex graph representations like disjunctive graphs [Zhang](#)
 2860

2861 ²⁰<https://github.com/henry-yeh/DeepACO>

et al. (2020), bipartite graphs Kwon et al. (2021), or heterogeneous graphs Song et al. (2022), making it harder to encode the problem. With RL4CO, we aim to mitigate these entry barriers for NCO researchers by providing established solution methods along with the environments. Further, by being modular by design, RL4CO allows for quick evaluation of different learning algorithms and network architectures, which can already lead to substantial improvements of the solution quality, as demonstrated in the example of the FJSSP in Table 3. Lastly, by providing benchmark instances like Taillard Taillard (1993) and easy ways of initializing the environments with external benchmark files, we facilitate the comparison of models with existing methods. The following chapter describes established DRL models for scheduling problems as well as their performance on synthetic and benchmark datasets.

E.3.1 JSSP

Models To solve the JSSP using DRL methods, we implement the L2D model described in Appendix C.2.7 in RL4CO. To train the encoder-decoder policy, we use the same Proximal Policy Optimization (PPO) algorithm as Zhang et al. (2020). In contrast to most other work in the NCO domain, L2D uses a (dense) stepwise reward function rather than a sparse episodic reward, which is observed only after a complete solution is obtained. This reward determines the change in the lower bound of the makespan given the partial schedule. Due to the dense nature of the reward, the PPO algorithm for the scheduling problems evaluates actions on a stepwise basis, whereas environments with an episodic reward are evaluated based on a full rollout. We compare these methods and discuss the different implementations in Appendix E.3.4.

Further, we demonstrate RL4CO’s ability to effortlessly implement a state-of-the-art solver for JSSP instances by exchanging the GCN encoder used by Zhang et al. (2020) with the MatNet encoder Kwon et al. (2021) described in Appendix C.2.11. Furthermore, the greedy decoding scheme of Zhang et al. (2020) is replaced by $N = 100$ random samples, of which the best is selected.

Reproduction and Improvement of Original Results We demonstrate RL4CO’s capability of learning dispatching rules for the JSSP by training and validating the L2D model of Zhang et al. (2020) and our version of L2D with the MatNet encoder on synthetic data. We report the performance achieved with RL4CO together with the baselines the authors of the original papers used, as well as the solutions obtained via the CP-Sat solver Google OR-Tools. The baselines are a set of selected PDRs that have a high practical relevance, namely Most Work Remaining (MWKR) and Most Operations Remaining (MOR).

Table 17: Comparison of RL4CO with L2D Zhang et al. (2020) and other baselines on the JSSP. For OR-Tools, the fraction of instances solved optimally is reported in parentheses.

Size	Metric	OR-Tools	PDRs		L2D	RL4CO	
			MWKR	MOR	Zhang et al. (2020)	GCN	MatNet ($\times 128$)
6×6	Obj.	487.75 (100%)	656.96	630.19	574.09	569.53	515.11
	Gap	-	34.6%	29.2%	17.7%	16.8%	5.6%
10×10	Obj.	808.32 (100%)	1151.41	1101.08	988.58	972.35	865.78
	Gap	-	42.6%	36.5%	22.3%	20.3%	7.1%
15×15	Obj.	1187.06 (99%)	1812.13	1693.33	1504.79	1492.94	1318.25
	Gap	-	52.6%	42.6%	26.7%	25.7%	11.0%
20×20	Obj.	1555.79 (4%)	2469.19	2263.68	2007.76	1992.36	1847.33
	Gap	-	58.6%	45.5%	29.0%	28.1%	18.7%

The results are listed in Table 17. RL4CO’s implementation of L2D manages to outperform the original implementation on all instance types, even when using the same model architecture, learning algorithm, and hyperparameters. The reason is that RL4CO uses an improved implementation of the environment. In the implementation of Zhang et al. (2020) the state of the environment does not contain a time dimension. Instead, the environment schedules the selected operation at the earliest feasible start time, given the current schedule. Here, we use the environment proposed by Tassel et al. (2021), where the environment transitions through distinct time steps $t = 0, 1, \dots, T$. In this

case, the start time of a selected operation is set to the time step at which it was selected, leading to a more natural form of credit assignment.

Using the MatNet encoder instead of the GCN and employing a decoding scheme based on multiple random rollouts further reduces the makespan by a large margin. One instances of size 6×6 , the gap to the optimal solutions was reduced by 11 percentage points to 5.6%, which corresponds to a third of the gap realized with the GCN encoder.

Taillard Benchmark and out-of-distribution performance With RL4CO, we also provide the possibility to test models against established benchmarks. For the JSSP, a well-recognized benchmark is that of Taillard (1993), which is also used by Zhang et al. (2020) to validate their model. In Table 18, we report the results of RL4CO on these instances along with the results obtained by Zhang et al. (2020) as well as the MOR and MWKR heuristics. We trained our MatNet models on JSSP instances up to size 20×20 . For larger Taillard instances, we report the out-of-distribution performance to demonstrate the model’s generalization ability. Similar to the synthetic test instances, our RL4CO implementation paired with the MatNet encoder manages to outperform the original L2D by large margins on all instances of the Taillard benchmark dataset, even when evaluating it on out-of-distribution instances.

Table 18: Results on the Taillard Taillard (1993) benchmark instances. BKS refers to the best known solutions and % opt. specifies the rate of instances with optimal solutions. Values marked with a \dagger indicate out-of-distribution performance of the model trained on 20×20 .

Size	Metric	BKS	PDRs		L2D	RL4CO
			MWKR	MOR	Zhang et al. (2020)	MatNet ($\times 128$)
15×15	Obj.	1230.06 (100%)	1927.5	1782.3	1547.50	1404.30
	Gap	-	56.7%	45.0%	26.0%	14.2%
20×15	Obj.	1363.22 (90%)	2190.7	2015.8	1774.7	1570.70
	Gap	-	60.7%	47.7%	30.0%	15.2%
20×20	Obj.	1617.60 (30%)	2518.6	2309.9	2128.1	1842.90
	Gap	-	55.7%	42.8%	31.6%	13.9%
30×15	Obj.	1787.68 (70%)	2728.0	2601.3	2378.8	2121.19 \dagger
	Gap	-	52.6%	45.6%	33.0%	18.6%
30×20	Obj.	1948.32 (0%)	3193.3	2888.1	2603.9	2357.90 \dagger
	Gap	-	63.9%	48.2%	33.6%	21.0%

E.3.2 FJSSP

Model To solve the FJSSP using DRL methods, we implement the HGNN model described in Appendix C.2.10 in RL4CO and train it with the same PPO algorithm as L2D. Besides HGNN we also implement a second model which exchanges the encoder of HGNN with the MatNet encoder.

Reproduction and Improvement of Original Results We compare the results obtained via RL4CO with those reported by Song et al. (2022) and the baseline used by them. Also, Song et al. (2022) use MWKR and MOR to benchmark their model as well as the OR-Tools solver. The results, which are obtained on a test set comprising of 100 randomly generated instances, are listed below in Table 19.

Similar to the JSSP, the HGNN implemented in RL4CO achieves better results than the original implementation, although both implementations use the same definition of the environment. However, in RL4CO, we use instance normalization Ulyanov et al. (2016) on the input variables as well as between consecutive HGNN layers, which we found to drastically stabilize the training process.

Again, we were able to enhance the quality of the solution further by simply exchanging the encoder with MatNet. Especially on the larger instances, the increased model complexity translates into much better model performance, with the solutions even surpassing OR-Tools on 20×10 instances.

2970 Table 19: Comparison of RL4CO and HGNN Song et al. (2022) on the FJSSP. For OR-Tools, the fraction
 2971 of instances solved optimally is reported in parentheses. Both RL4CO and Song et al. (2022) make use of
 2972 random-rollouts for decoding.

2973

Size	Metric	OR-Tools	PDRs		HGNN	RL4CO ($\times 128$)	
			MWKR	MOR	Song et al. (2022) ($\times 128$)	HGNN	MatNet
10×5	Obj.	96.59 (15%)	115.29	116.69	105.61	102.49	99.02
	Gap	-	19.4%	20.9%	9.4%	6.1%	2.5%
20×5	Obj.	188.45 (0%)	216.98	217.17	207.50	199.47	192.05
	Gap	-	15.2%	15.3%	10.1%	5.8%	1.9%
15×10	Obj.	145.42 (5%)	169.18	173.40	160.36	155.34	151.93
	Gap	-	16.3%	19.3%	10.3%	6.8%	4.5%
20×10	Obj.	197.24 (0%)	220.85	221.86	214.87	207.52	192.00
	Gap	-	11.9%	12.53%	9.0%	5.2%	-2.7%

2985

2986 **Out-of-distribution** In this section, we evaluate the out-of-distribution performance of the DRL
 2987 models trained with RL4CO on FJSSP 20×10 instances, by evaluating them on smaller (20×5
 2988 & 15×10) and larger (30×10 & 40×10) instances. The results in Table 20 indicate that both
 2989 HGNN and MatNet manage to generalize well to problems of different sizes. Despite being trained
 2990 on smaller instances, the HGNN manages to close the performance gap when evaluated on larger
 2991 instances, with gaps being as small as 3.7% for FJSSP 40×10 instances. And on FJSSP 20×5
 2992 instances, the average makespan increases by only 1.56 (0.8%) when using the model trained on
 2993 FJSSP 20×10 instead of 20×5 instances. Again, the MatNet model shows superior perform-
 2994 ance compared to the other baselines and surpasses even the results obtained by OR-Tools on the
 2995 larger instances. The within-distribution performance of MatNet, therefore, also translates to out-of-
 2996 distribution instances, indicating that the complexity of the model results in a better generalization
 2997 ability.

2998 Table 20: Generalization performance of a policy trained on a 20×10 FJSSP instances on smaller and larger
 2999 instances. We use 100 test instances per instance size. Gaps are reported with respect to the results of OR-Tools

3000

Size	Metric	OR-Tools	PDRs		HGNN	RL4CO ($\times 128$)	
			MWKR	MOR	Song et al. (2022) ($\times 128$)	HGNN	MatNet
20×5	Obj.	188.45 (0%)	216.98	217.17	207.50	201.03	193.61
	Gap	-	15.2%	15.3%	10.1%	6.7%	2.7%
15×10	Obj.	145.42 (5%)	169.18	173.40	160.36	162.41	150.59
	Gap	-	16.3%	19.3%	10.3%	11.7%	3.5%
30×10	Obj.	294.10 (0%)	319.89	320.18	312.20	309.10	286.16
	Gap	-	8.8%	8.9%	6.1%	5.1%	-2.7%
40×10	Obj.	397.36 (0%)	425.70	425.19	415.14	412.05	381.19
	Gap	-	7.1%	7.0%	4.4%	3.7%	-4.1%

3012 E.3.3 FFSP

3013

3014

3015 **MatNet** To solve the FFSP using DRL, RL4CO implements the policy network described by
 3016 Kwon et al. (2021). It uses separate policy networks for each stage of the FFSP. Each of the stage
 3017 networks employs the MatNet encoder described in Appendix C.2.11, which generates embeddings
 3018 for jobs and machines using the processing times of the job-machine pairs of the respective stage.
 3019 The decoder of the attention model Kool et al. (2019a) then utilizes the machine embeddings of
 3020 the respective stage as query and the job embeddings as keys and values to compute the probability
 3021 distribution over jobs.

3022

3023 **Results** We use the same three instance types described by Kwon et al. (2021) to evaluate our
 implementations of the FFSP environment and the policy network. The instances only differ in

the number of jobs, which are set to 20, 50, and 100. We assume that there are $S = 3$ stages, and each stage has $M = 4$ machines. In the k th stage, the processing time of the job j on the machine m is given by p_{jmk} . Therefore, an instance of the problem is defined by three matrices (P_1 , P_2 , and P_3), specifying the processing time for each job-machine combination in that stage. We report the results obtained by RL4CO and compare them to those obtained by Kwon et al. (2021) in Table 21. Other benchmarks used are the exact solver CPLEX (for which results can only be obtained for FFSP20 instances), the Shortest Job First (SJF) dispatching rule, as well as the evolutionary algorithms Particle Swarm Optimization (PSO), and Genetic Algorithm (GA). One can see that, using RL4CO, we are able to reproduce the results from the original paper.

Table 21: Comparison of RL4CO with the results reported in Kwon et al. (2021). Gaps are reported with respect to the best known results.

Instance	Matric	CPLEX (600s)	SJF	GA	PSO	Kwon et al. (2021)	RL4CO
FFSP20	Obj.	36.6	31.3	30.6	29.1	27.3	27.2
	Gap	34.5%	15.0%	12.5%	6.9%	0.3%	0.0%
FFSP50	Obj.	-	57.0	56.4	55.1	51.5	51.6
	Gap	-	10.7%	9.5%	7.0%	0.0%	0.2%
FFSP100	Obj.	-	99.3	98.7	97.3	91.5	91.3
	Gap	-	8.8%	8.1%	6.6%	0.2%	0.0%

E.3.4 DENSE AND EPISODIC REWARDS

We additionally compare dense and episodic rewards for the TSP and FJSSP environments, with similar training settings as in other experiments, except for the different reward functions.

Here, we compare the performance of the HGNN (Song et al., 2022) in solving the FJSSP and AM (Kool et al., 2019a) in solving the TSP when trained using a stepwise vs. an episodic reward. The results in Table 22 show that evaluating the FJSSP in a stepwise manner and stepwise re-encoding the current state significantly outperforms a policy based on a single, episodic reward. This is reasonable since the state of the FJSSP has many dynamic elements, and a policy that relies on a single encoder step may not fully grasp the problem dynamics. On the other hand, stepwise rewards for the TSP (AM model trained with POMO with the settings as Kwon et al. (2020)) do not work well, and interestingly, performance approaches roughly that of the nearest insertion algorithms. Different CO problems react to the same learning setup, which again underpins the importance of a unified framework where different algorithms are implemented and are easily exchangeable.

Table 22: Comparison of dense (i.e. stepwise) and episodic rewards for the TSP and the FJSSP

Reward	TSP			FJSSP		
	20	50	100	10×5	20×5	15×10
Dense	4.51	7.05	9.80	102.49	199.47	155.34
Episodic	3.83	5.81	7.82	110.65	204.88	182.90

E.4 ELECTRONIC DESIGN AUTOMATION: LEARNING TO PLACE DECAPS

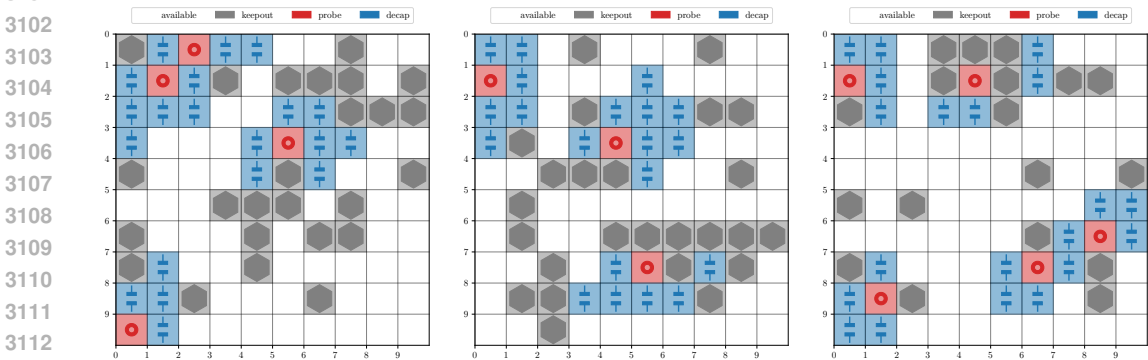
Setup In this section, we benchmark models on the mDPP from Appendix B.3.2. We benchmark 3 variants of online DevFormer (DF), namely DF(PG,Critic): REINFORCE (where PG stands for Policy Gradients, an “alias” of the REINFORCE algorithm) with Critic baseline, DF(PG,Rollout): REINFORCE with Rollout baseline as well as PPO. All experiments are run with the same hyperparameters as the other experiments except for the batch size set to 64, the maximum number of samples set to 10,000, and a total of only 10 epochs due to the nature of the benchmark sample efficiency.

3078 E.4.1 MAIN RESULTS
3079

3080 **Table 23** shows the main numerical results for the task when RS, GA, and DF models are trained for
3081 placing 20 decaps. While RS and GA need to take online shots to solve the problems (we restricted
3082 the number to 100), DF models can successfully predict in a zero-shot manner and outperform the
3083 classical approaches. Interestingly, the vanilla critic-based method performed the worst, while our
3084 implementation of PPO almost matched the rollout policy gradients (PG) baseline; since extensive
3085 hyperparameter tuning was not performed, we expect PPO could outperform the rollout baseline
3086 given it requires fewer samples. **Fig. 24** shows example renderings of the solved environment.

3087
3088 Table 23: Performance of different methods on the mDPP benchmark
3089

Method	# Shots	Score \uparrow	
		maxsum	maxmin
<i>Online Test Time Search</i>			
Random Search	100	11.55	10.63
Genetic Algorithm	100	11.93	11.07
<i>RL Pretraining & Zero Shot Inference</i>			
DF-(PG,Critic)	0	10.89 \pm 0.63	9.51 \pm 0.68
DF-(PPO)	0	12.16 \pm 0.03	11.17 \pm 0.11
DF-(PG,Rollout)	0	12.21 \pm 0.01	11.26 \pm 0.03



3113
3114 **Figure 24:** Renders of the environment with *maxmin* objective solved by DF-(PG,Rollout). The model suc-
3115 cessfully learned one main heuristic for DPP problems, which is that the optimal placement of decaps (blue) is
3116 generally close to probing ports (red).

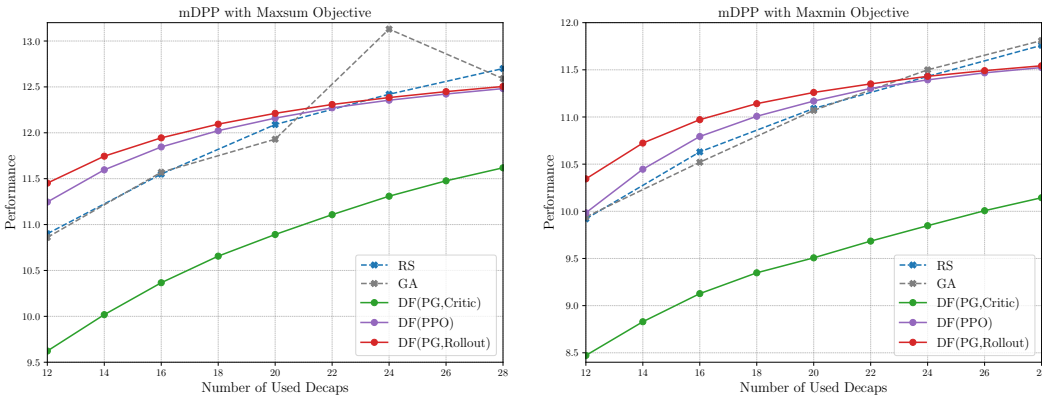
3117
3118
3119 E.4.2 GENERALIZATION TO DIFFERENT NUMBER OF COMPONENTS
3120

3121 In hardware design, the number of components is one major contribution to cost; ideally, one would
3122 want to use the least number of components possible with the best performance. In the DPP, in-
3123 creasing the number of decaps *generally* improves the performance at a greater cost, hence Pareto-
3124 efficient models are essential to identify. **Fig. 25** shows the performance of DF models trained on
3125 20 decaps against the baselines. DF models PPO and PG-rollout can successfully generalize and are
3126 also Pareto-efficient with fewer decaps, important in practice for cost and material saving.

3127
3128 E.5 LEARNING TO IMPROVE
3129

3130 In this section, we first show the efficiency of RL4CO when reproducing the improvement methods
3131 on the TSP and PDP with 50 nodes and discuss the potential collaboration of constructive methods
with improvement methods for better inference performance.

3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144



3145 Figure 25: Performance vs number of used decaps for mDPP with *maxsum* objective [Left] and *maxmin* objective [Right].

3146
3147
3148
3149 E.5.1 MAIN RESULTS

3150
3151 As shown in Table 24, refactoring and implementing the three improvement methods—DACT Ma et al. (2021) (TSP50), N2S Ma et al. (2022) (PDP50), and NeuOpt Ma et al. (2024) (PDP50)—using RL4CO consistently results in better efficiency compared to the original implementations. Specifically, training and testing times ($T = 1,000$) are faster, and peak memory usage is lower. This advancement can be attributed to RL4CO’s streamlined design, which uses a single tensor dictionary variable to store all state information, and the incorporation of efficient libraries like PyTorch Lightning and TorchRL. These enhancements demonstrate RL4CO’s superior efficiency and ease of implementation.

3152
3153
3154
3155
3156
3157
3158
3159 Table 24: Comparison of time and memory usage for DACT Ma et al. (2021) (TSP50), N2S Ma et al. (2022) (PDP50), and NeuOpt Ma et al. (2024) (PDP50) between the original implementation and the RL4CO implementation.

3160
3161
3162
3163

	T_train (one epoch)	T_test (1k,1k)	Memory
DACT-Origin	16m	38s	8069MB
DACT-RL4CO	10m	26s	7135MB
N2S-Origin	26m	41s	13453MB
N2S-RL4CO	17m	33s	12489MB
NeuOpt-Origin	14m	37s	7273MB
NeuOpt-RL4CO	10m	31s	6313MB

3164
3165
3166
3167
3168
3169
3170

3171
3172 E.5.2 DISCUSSION

3173
3174 As shown in Fig. 26, bootstrapping improvement with constructive methods can greatly improve the performance, especially in terms of the Primal Integral (PI, Appendix D.1.2). While in TSP bootstrapping is consistently better than simply improving with default solutions (i.e. lower final gap to BKS as well as PI), we note that in PDP with N2S, improving starting from a random initialization can yield better performance in terms of gap. However, the PI reveals that while N2S from random init achieves a value of 5.580, N2S from HAM construction initialization achieves a much better 2.234, indicating a much better early convergence speed and Pareto front.

3175
3176
3177
3178
3179
3180
3181 We additionally offer some clues on how to improve such performance. Firstly, we simply initialized from a greedy solution, while more complex inference strategies may offer a significant boost. Furthermore, the trained model as per the setting in Appendix D.3.3 could be further trained and obtain better performance. Importantly, we believe that *end-to-end construction & improvement*, in which both a constructive and improvement method are trained together, could ultimately outperform a separate training and achieve the best of both worlds.

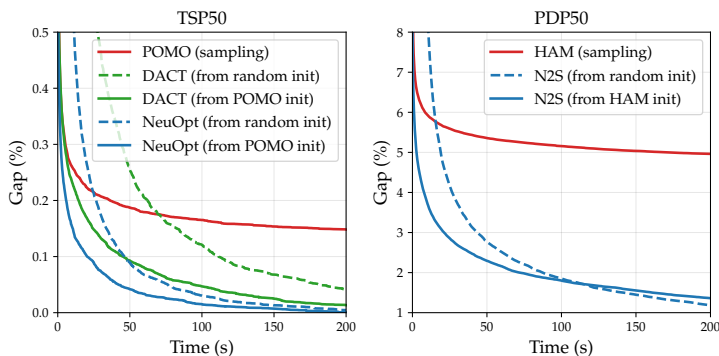


Figure 26: Bootstrapping improvement with constructive methods for TSP50 and PDP50.

E.6 GRAPH PROBLEMS: FACILITY LOCATION PROBLEM (FLP) AND MAXIMUM COVERAGE PROBLEM (MCP)

Here, we present the experimental results and the corresponding discussions on the two CO problems on graphs: the Facility Location Problem (FLP; see [Appendix B.4.1](#)) and the Maximum Coverage Problem (MCP; see [Appendix B.4.2](#)).

E.6.1 EXPERIMENTAL SETTINGS

Baseline methods We consider two simple baselines: uniform random (UR) and deterministic greedy (DG), where UR chooses k locations uniformly at random and DG chooses k locations one by one in a greedy manner. We also apply two MIP solvers, Gurobi ([Gurobi Optimization, 2021](#)) and SCIP ([Bestuzheva et al., 2021](#)), to obtain the optimal solutions.

Benchmark methods We benchmark with the attention model (AM) with different embedding models (i.e., encoders) and different RL baselines. For FLP, the considered embedding models are: the multilayer perceptron (MLP), the graph convolutional network (GCN) ([Kipf and Welling, 2017](#)), and the graph attention network ([Velickovic et al., 2017](#); [Brody et al., 2019](#)). For MCP, since the problem instances are formulated on bipartite graphs, the considered embedding models are: the multilayer perceptron (MLP), the GraphSAGE model ([Hamilton et al., 2017](#)) (in short “SAGE”), and the generalized GCN model ([Li et al., 2020a](#)) (in short “GEN”). The considered RL baselines are: Rollout, Mean, Exponential, and Critic. All the models are trained in 100 epochs. The learning rate is $1e-5$ for FLP and $1e-4$ for MCP. In each epoch, 100,000 training data are used with batch size 1,000. For the decoding strategies, we consider sampling (with 64 independent samples) and greedy. For sampling (and UR), we report both the “best” performance among the 64 independent samples and the “mean” (i.e., average) performance over the 64 independent samples.

Test-time active search We apply three variants of active search at test time: the original active search (AS) proposed by [Bello et al. \(2017\)](#), efficient active search (EAS) proposed by [Hottung et al. \(2022\)](#) with two variants: EAS-Emb that finetunes embeddings and EAS-Lay that finetunes new layers. We run all the active search variants for 100 iterations.

E.6.2 BENCHMARK RESULTS

Main benchmark [Table 25](#) shows the main numerical results when the methods are trained and tested to choose $k = 10$ locations on instances with $n = 100$ locations. [Table 26](#) shows the main numerical results when the methods are trained and tested to choose $k = 10$ sets on instances with $n = 100$ sets and $m = 200$ items in total. Each item has a random weight between 1 and 10, and the number of items in each set is randomly sampled between 5 and 15. The reported results are averaged over 1,000 randomly generated test instances. We also report the average gap between the performance for each setting and the optimum by solvers as described in [Appendix D.1.1](#).

Here we use absolute values since we *minimize* the total distance for FLP while *maximizing* the total weights for MCP. When using absolute values, it is consistent that smaller gaps correspond to

3240 Table 25: Performance of different methods on the facility location problem (FLP) benchmark. For the perfor-
 3241 mance, the smaller the better.

Encoder	RL Baseline	Sample (Best)	Sample (Mean)	Greedy	AS	Active Search	
						EAS-Emb	EAS-Lay
MLP	Rollout	10.4895	11.0056	10.9980	10.3004	10.2997	10.2997
	(Gap)	(2.19%)	(7.23%)	(7.16%)	(0.35%)	(0.34%)	(0.34%)
	Mean	10.5635	11.1614	10.9350	10.2995	10.3008	10.3008
	(Gap)	(2.91%)	(8.75%)	(6.54%)	(0.34%)	(0.35%)	(0.35%)
	Exponential	10.5726	11.1848	10.9589	10.3054	10.3051	10.3051
	(Gap)	(3.00%)	(8.98%)	(6.78%)	(0.40%)	(0.39%)	(0.39%)
GCN	Critic	10.5617	11.1401	10.9439	10.2987	10.2994	10.2994
	(Gap)	(2.90%)	(8.55%)	(6.63%)	(0.33%)	(0.34%)	(0.34%)
	Rollout	10.4232	10.6404	10.6094	10.2955	10.2956	10.2958
	(Gap)	(1.54%)	(3.66%)	(3.36%)	(0.30%)	(0.30%)	(0.30%)
	Mean	10.4321	10.8095	10.6076	10.2807	10.2830	10.2830
	(Gap)	(1.63%)	(5.31%)	(3.34%)	(0.15%)	(0.18%)	(0.18%)
GAT	Exponential	10.4729	10.9573	10.7257	10.2837	10.2859	10.2859
	(Gap)	(2.02%)	(6.75%)	(4.49%)	(0.18%)	(0.20%)	(0.20%)
	Critic	10.7086	11.4549	11.0139	10.2859	10.2891	10.2891
	(Gap)	(3.82%)	(0.54%)	(6.01%)	(0.20%)	(0.23%)	(0.23%)
	Rollout	10.4685	10.9202	10.8916	10.2956	10.2956	10.2957
	(Gap)	(1.99%)	(6.40%)	(6.12%)	(0.30%)	(0.30%)	(0.30%)
GAT	Mean	10.6641	11.3499	11.0133	10.2865	10.2899	10.2898
	(Gap)	(3.90%)	(0.59%)	(7.31%)	(0.21%)	(0.24%)	(0.24%)
	Exponential	10.6487	11.3504	10.9869	10.2864	10.2881	10.2880
	(Gap)	(3.75%)	(0.60%)	(7.05%)	(0.21%)	(0.22%)	(0.22%)
	Critic	10.6566	11.3440	10.8813	10.2859	10.2888	10.2888
	(Gap)	(4.33%)	(1.62%)	(7.31%)	(0.20%)	(0.23%)	(0.23%)
Uniform Random (Best)						12.4788	
(Gap)						(21.62%)	
Uniform Random (Mean)						15.6327	
(Gap)						(52.40%)	
Deterministic Greedy						10.9831	
(Gap)						(7.02%)	
GUROBI/SCIP (Optimum)						10.2650	
(Gap)						(0.00%)	

3273
 3274
 3275 better performance. The performance of RL methods with sampling is consistently better than the
 3276 two baselines, uniform random (UR) and deterministic greedy (DG), showing their effectiveness on
 3277 those two problems.

3278
 3279 **Effect of the encoder** Overall, the performance of different encoders is similar. For FLP, we can
 3280 observe GCN’s marginal superiority (except when we use Critic as the RL baseline). For MCP, the
 3281 best encoders for different RL baselines are different, but MLP’s performance is the overall best.

3282
 3283 **Effect of the RL baseline** For FLP, for the four considered RL baselines (Rollout, Mean, Expo-
 3284 nential, Critic), Rollout is consistently better than the other three. For MCP, the differences in the
 3285 performance of different RL baselines are not significant.

3286
 3287 **Effect of active search** Active search significantly improves performance in almost all cases. For
 3288 FLP, interestingly, Rollout achieves the best overall performance without active search, but Rollout
 3289 underforms in many cases with test-time active search. Notably, the performance of the original
 3290 active search (AS) is less stable than the two variants of efficient active search (EAS), especially for
 3291 MCP. In our understanding, AS was originally designed for routing problems and uses multi-start
 3292 sampling with distinct initial action (i.e., the first location/set to choose). Such a strategy is useful
 3293 for routing problems due to symmetry but is less useful for problems without symmetry, such as
 FLP and MCP.

Table 26: Performance of different methods on the maximum coverage problem (MCP) benchmark. For the performance, the larger the better.

Encoder	RL Baseline	Sample (Best)	Sample (Mean)	Greedy	AS	Active Search EAS-Emb	EAS-Lay
MLP	Rollout	682.4741	662.4359	665.1740	689.6200	689.6070	689.6070
	(Gap)	(0.96%)	(3.31%)	(3.05%)	(0.09%)	(0.09%)	(0.09%)
	Mean	682.4011	664.7105	668.7470	682.0610	689.5900	689.5900
	(Gap)	(1.06%)	(3.96%)	(3.56%)	(1.18%)	(0.09%)	(0.09%)
	Exponential	683.0300	665.1467	666.6640	671.3130	689.5870	689.5870
	(Gap)	(1.09%)	(3.99%)	(3.64%)	(9.68%)	(0.09%)	(0.09%)
	Critic	683.1511	666.9047	668.6411	687.8240	689.3510	689.3510
	(Gap)	(1.43%)	(5.40%)	(4.92%)	(0.35%)	(0.13%)	(0.13%)
	SAGE	Rollout	681.8690	664.1233	665.9901	689.4810	689.5020
(Gap)		(1.14%)	(3.71%)	(3.44%)	(0.11%)	(0.11%)	(0.11%)
Mean		682.1360	669.2791	670.4091	666.0360	689.5990	689.5890
(Gap)		(1.06%)	(3.63%)	(3.05%)	(10.44%)	(0.09%)	(0.09%)
Exponential		680.3970	653.0383	656.3170	675.2220	689.5990	689.5980
(Gap)		(1.06%)	(3.95%)	(3.46%)	(2.18%)	(0.09%)	(0.09%)
Critic		676.9190	645.9108	649.6940	647.9050	688.4500	688.4650
(Gap)		(1.94%)	(6.43%)	(5.89%)	(6.12%)	(0.26%)	(0.26%)
GEN		Rollout	680.2640	648.2318	656.3710	689.4430	689.4660
	(Gap)	(1.10%)	(2.96%)	(2.80%)	(0.12%)	(0.11%)	(0.11%)
	Mean	682.1960	662.1896	664.6721	681.3950	689.5670	689.5670
	(Gap)	(0.97%)	(3.56%)	(3.34%)	(1.28%)	(0.10%)	(0.10%)
	Exponential	682.4290	662.5012	665.8010	689.4060	689.5650	689.5650
	(Gap)	(1.07%)	(3.70%)	(3.18%)	(0.12%)	(0.10%)	(0.10%)
	Critic	682.3510	664.1604	667.7340	689.6170	689.3940	689.3940
	(Gap)	(1.45%)	(6.08%)	(4.91%)	(0.09%)	(0.12%)	(0.12%)
	Uniform Random (Best)					527.9360	
(Gap)					(-23.50%)		
Uniform Random (Mean)					432.7287		
(Gap)					(-37.30%)		
Deterministic Greedy					680.2050		
(Gap)					(-1.46%)		
GUROBI/SCIP (Optimum)					690.2350		
(Gap)					(0.00%)		

Test-time sampling techniques We also consider other test-time sampling techniques: top- p sampling (Holtzman et al., 2019) and different sampling temperatures. Top- p sampling discards actions with low probabilities, and top- p sampling with lower p values discards more low-probability actions. For sampling temperatures, higher temperatures give more uniform sampling. The considered p values are: 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.99, 1.0. The sampling temperatures considered are 0.01, 0.03, 0.1, 0.3, 0.5, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.5, and 2.0. Fig. 27 show the heatmaps for each combination of encoder and RL baseline, for FLP and MCP. In each subplot, the x -axis represents the value of p in top- p sampling, and the y -axis represents the sampling temperature. For each combination, the best performance is marked with a red star. For FLP, the best performance is usually achieved with a proper (i.e., neither too high nor too low) level of randomness. As the p value of top- p sampling increases, the best sampling temperature decreases. Recall that both increasing the p value and increasing the sampling temperature would increase the randomness in sampling. Overall, compared to other RL baselines, Rollout needs a higher level of randomness to perform best. For MCP, the best performance is usually achieved without top- p sampling and with a high sampling temperature, i.e., without high randomness in the sampling space.

E.6.3 OUT-OF-DISTRIBUTION

Results on out-of-distribution instances Table 27 shows the main numerical results when the methods are trained to choose $k = 10$ locations on instances with $n = 100$ locations, but tested to choose $k' = 20$ locations on instances with $n' = 200$ locations. Table 28 shows the main numerical results when the methods are trained to choose $k = 10$ sets on instances with $n = 100$ sets and

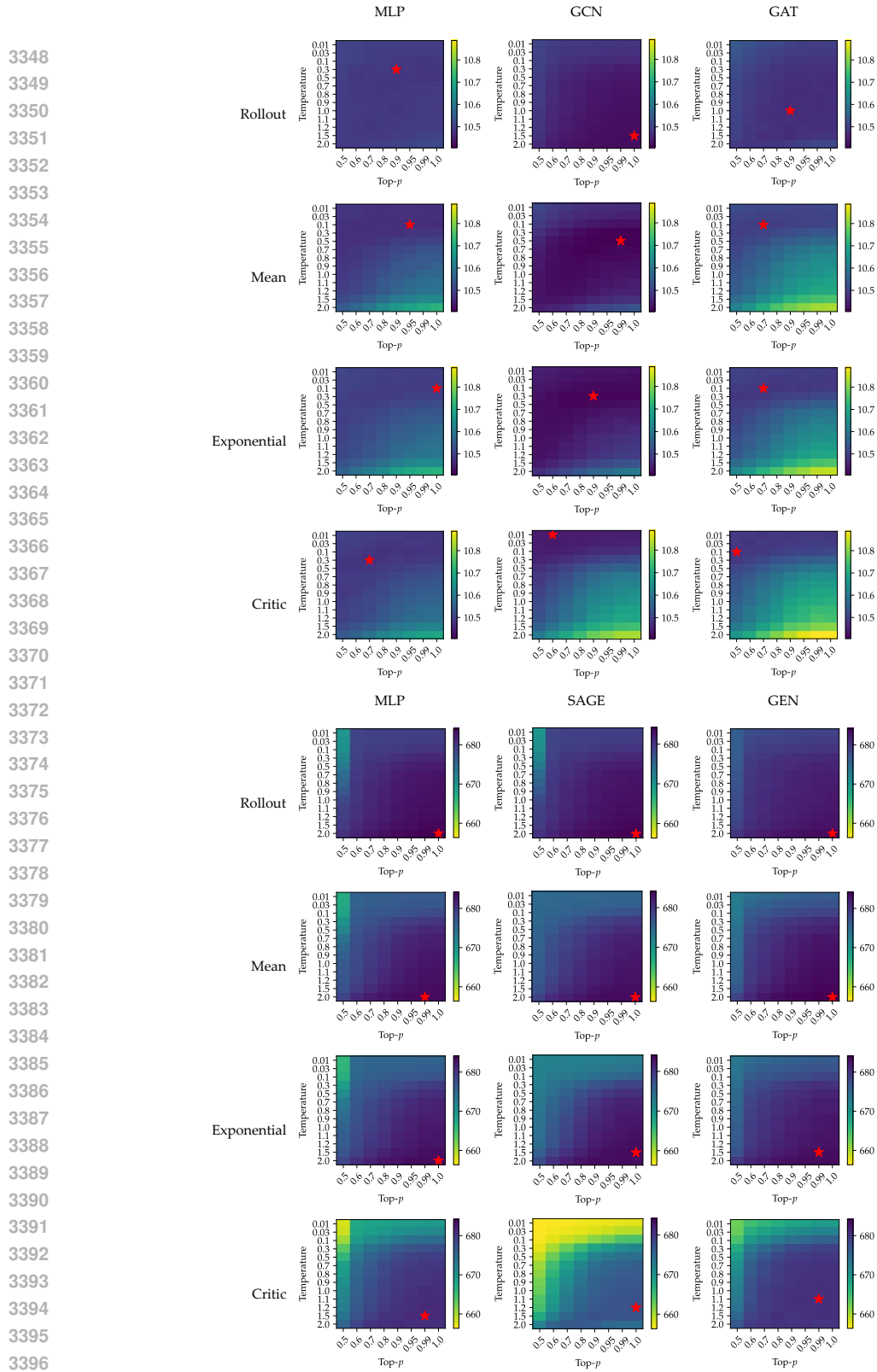


Figure 27: Performance of sampling with different p values for top- p sampling and different sampling temperatures. Top: FLP; Bottom: MCP. For each combination of encoder and RL baseline, the best performance is marked with a star.

Table 27: Performance of different methods on the facility location problem (FLP) out-of-distribution instances. For the performance, the smaller the better.

Encoder	RL Baseline	Sample (Best)	Sample (Mean)	Greedy	Active Search			
					AS	EAS-Emb	EAS-Lay	
MLP	Rollout	14.7612	15.2979	15.2709	14.4160	14.4181	14.4181	
	(Gap)	(3.85%)	(7.63%)	(7.44%)	(1.42%)	(1.43%)	(1.43%)	
	Mean	15.0045	15.7343	15.3075	14.5315	14.5331	14.5331	
	(Gap)	(5.56%)	(10.70%)	(7.70%)	(2.23%)	(2.24%)	(2.24%)	
	Exponential	15.0022	15.7144	15.3131	14.5274	14.5266	14.5266	
	(Gap)	(5.54%)	(10.56%)	(7.74%)	(2.20%)	(2.19%)	(2.19%)	
	Critic	14.9670	15.6631	15.2781	14.5147	14.5132	14.5132	
	(Gap)	(5.30%)	(10.20%)	(7.49%)	(2.11%)	(2.10%)	(2.10%)	
	GCN	Rollout	14.9564	15.4230	15.3610	14.6254	14.6239	14.6248
		(Gap)	(5.22%)	(8.51%)	(8.07%)	(2.89%)	(2.88%)	(2.89%)
		Mean	15.1380	15.8310	15.3713	14.6554	14.6572	14.6574
		(Gap)	(6.50%)	(11.38%)	(8.14%)	(3.10%)	(3.11%)	(3.12%)
Exponential		15.2197	15.9598	15.4441	14.6961	14.6963	14.6973	
(Gap)		(7.08%)	(12.29%)	(8.66%)	(3.39%)	(3.39%)	(3.40%)	
Critic		15.1754	15.9835	15.2815	14.6579	14.6634	14.6642	
(Gap)		(6.53%)	(12.00%)	(8.23%)	(3.12%)	(3.16%)	(3.16%)	
GAT		Rollout	14.7503	15.2808	15.2593	14.4142	14.4150	14.4143
		(Gap)	(3.77%)	(7.51%)	(7.36%)	(1.40%)	(1.41%)	(1.40%)
		Mean	15.1147	15.9092	15.2895	14.5944	14.5986	14.5946
		(Gap)	(6.34%)	(11.93%)	(7.57%)	(2.67%)	(2.70%)	(2.67%)
	Exponential	15.1639	15.9886	15.2945	14.5991	14.6004	14.6011	
	(Gap)	(6.68%)	(12.49%)	(7.60%)	(2.70%)	(2.71%)	(2.72%)	
	Critic	15.1428	15.9191	15.3835	14.6053	14.6111	14.6111	
	(Gap)	(6.76%)	(12.46%)	(7.51%)	(2.75%)	(2.79%)	(2.79%)	
	Uniform Random (Best)					18.3215		
	(Gap)					(28.92%)		
	Uniform Random (Mean)					21.7044		
	(Gap)					(52.74%)		
Deterministic Greedy					15.3090			
(Gap)					(7.71%)			
GUROBI/SCIP (Optimum)					14.2148			
(Gap)					(0.00%)			

$m = 200$ items in total and tested to choose $k' = 20$ sets on instances with $n' = 200$ sets and $m' = 400$ items in total. Each item has a random weight between 1 and 10, and the number of items in each set is randomly sampled between 5 and 15. The reported results are averaged over 1,000 randomly generated test instances. We also report the average gap for each setting. Overall, the performance of RL methods generalizes well to out-of-distribution instances, being significantly higher than both Uniform Random and Deterministic Greedy with enough sampling.

Effect of the encoder For FLP, unlike the main benchmark, the superiority of GCN no longer exists for out-of-distribution instances. For MCP, the best encoders for different RL baselines are still different, and the performance of MLP is the best.

Effect of the RL baseline For FLP, again, Rollout is overall better than the other three. For MCP, the best RL baselines for different encoders are different, and Mean and Critic are overall good choices.

Effect of active search Again, active search clearly improves performance in almost all cases. For FLP, unlike the main benchmark, for out-of-distribution instances, Rollout overall performs best with and without active search. Still, the performance of the original active search (AS) is less stable than the two variants of efficient active search (EAS). With active search (especially EAS), the performance of RL methods is consistently better than that of Deterministic Greedy and is close to the optimum.

Table 28: Performance of different methods on the maximum coverage problem (MCP) out-of-distribution instances. For the performance, the larger the better.

Encoder	RL Baseline	Sample (Best)	Sample (Mean)	Greedy	AS	Active Search		
						EAS-Emb	EAS-Lay	
MLP	Rollout	1356.8970	1299.8690	1307.5250	1385.3340	1385.3280	1385.3280	
	(Gap)	(-1.83%)	(-5.48%)	(-5.03%)	(-0.32%)	(-0.33%)	(-0.33%)	
	Mean	1360.7710	1306.4015	1312.6290	1319.8180	1383.3580	1383.3580	
	(Gap)	(-2.34%)	(-6.45%)	(-5.89%)	(-5.04%)	(-0.47%)	(-0.47%)	
	Exponential	1360.7830	1306.3337	1312.7070	1088.0180	1383.9670	1383.9670	
	(Gap)	(-2.49%)	(-6.64%)	(-6.23%)	(-21.71%)	(-0.42%)	(-0.42%)	
	Critic	1363.9190	1313.2830	1319.5280	1353.9080	1377.3780	1377.3780	
	(Gap)	(-3.29%)	(-7.83%)	(-7.33%)	(-2.59%)	(-0.90%)	(-0.90%)	
	SAGE	Rollout	1353.9790	1297.5763	1303.7120	1382.2220	1382.1140	1382.1140
		(Gap)	(-2.55%)	(-6.61%)	(-6.16%)	(-0.55%)	(-0.56%)	(-0.56%)
Mean		1366.0050	1320.5641	1325.5570	1121.7650	1384.3780	1384.3650	
(Gap)		(-2.06%)	(-5.98%)	(-5.53%)	(-19.30%)	(-0.39%)	(-0.40%)	
Exponential		1344.1420	1281.0377	1288.0360	1288.2830	1383.6030	1383.5500	
(Gap)		(-2.30%)	(-6.38%)	(-5.73%)	(-7.31%)	(-0.45%)	(-0.45%)	
Critic		1331.1100	1266.6130	1276.0670	1092.0550	1367.4660	1367.4690	
(Gap)		(-4.23%)	(-8.87%)	(-8.19%)	(-21.42%)	(-1.61%)	(-1.61%)	
GEN		Rollout	1334.2700	1269.0966	1284.4550	1385.6540	1385.5750	1385.5750
		(Gap)	(-1.68%)	(-4.96%)	(-4.60%)	(-0.30%)	(-0.31%)	(-0.31%)
	Mean	1354.8450	1297.2153	1302.8560	1305.4070	1384.3080	1384.2980	
	(Gap)	(-2.06%)	(-5.98%)	(-5.52%)	(-6.08%)	(-0.40%)	(-0.40%)	
	Exponential	1357.4750	1300.7056	1309.8040	1376.1300	1384.3780	1384.3900	
	(Gap)	(-2.11%)	(-6.18%)	(-5.45%)	(-0.99%)	(-0.39%)	(-0.39%)	
	Critic	1360.0420	1303.4360	1313.6640	1366.2960	1374.8630	1374.8370	
	(Gap)	(-4.00%)	(-8.68%)	(-7.58%)	(-1.69%)	(-1.08%)	(-1.08%)	
	Uniform Random (Best)					1003.3390		
	(Gap)					(-27.80%)		
Uniform Random (Mean)					866.3536			
(Gap)					(-37.66%)			
Deterministic Greedy					1367.2240			
(Gap)					(-1.63%)			
GUROBI/SCIP (Optimum)					1389.8450			
(Gap)					(0.00%)			

Test-time sampling techniques For out-of-distribution instances, we also consider top- p sampling and different sampling temperatures as the main benchmark. The considered p values are: 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.99, 1.0. The sampling temperatures considered are 0.01, 0.03, 0.1, 0.3, 0.5, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.5, and 2.0. Fig. 28 show the heatmaps for each combination of encoder and RL baseline, for FLP and MCP. In each subplot, the x -axis represents the value of p in top- p sampling, and the y -axis represents the sampling temperature. For each combination, the best performance is marked with a red star. For both FLP and MCP, the best performance is usually achieved with a proper (i.e., neither too high nor too low) level of randomness. As the p value of top- p sampling increases, the best sampling temperature decreases. Recall that both increasing the p value and increasing the sampling temperature would increase the randomness in sampling.

E.7 EFFICIENT SOFTWARE ROUTINES

E.7.1 MIXED-PRECISION TRAINING

RL4CO supports multiple device types as well as floating point precisions by leveraging PyTorch Lightning (Falcon and The PyTorch Lightning team, 2019).

As Table 29 shows mixed-precision training can successfully reduce computational costs both in terms of runtime and especially with memory usage.

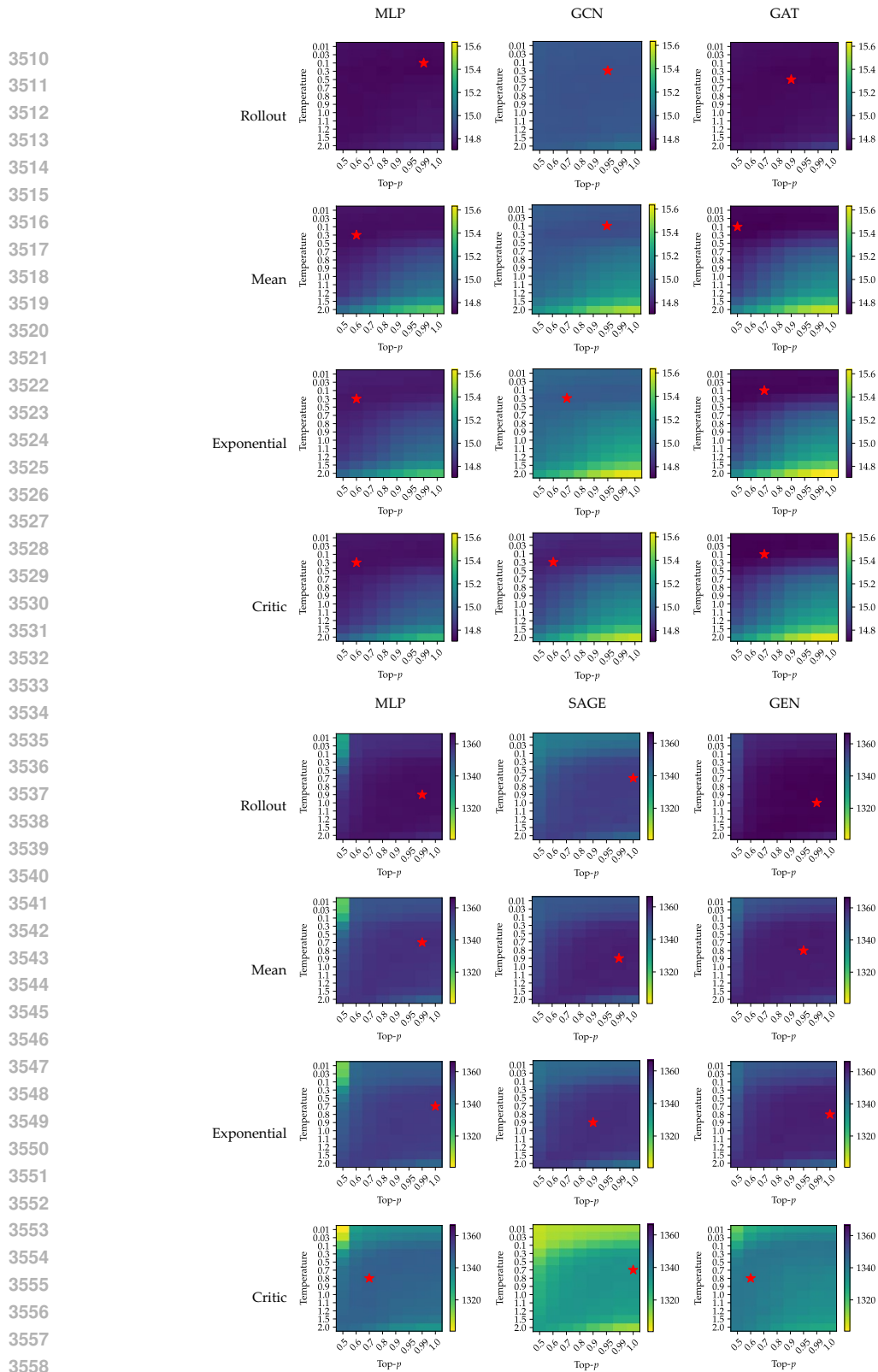


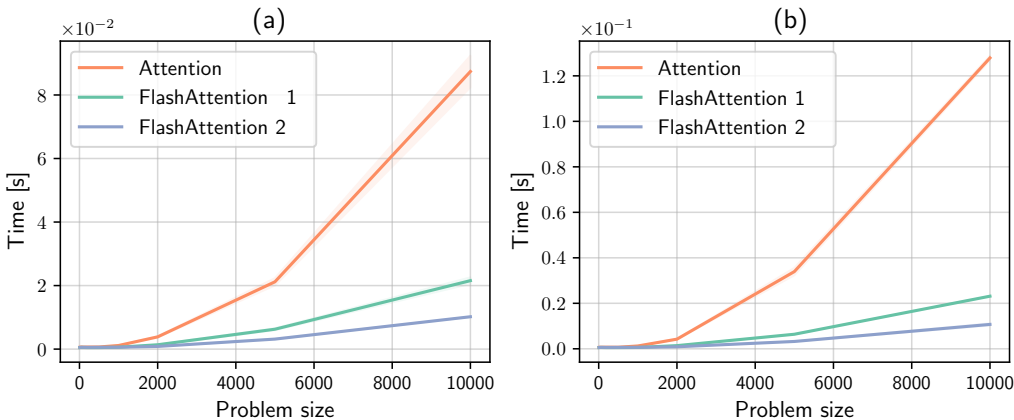
Figure 28: Performance of sampling on out-of-distribution instances with different p values for top- p sampling and different sampling temperatures. Top: FLP; Bottom: MCP. For each combination of encoder and RL baseline, the best performance is marked with a star.

3564 Table 29: Running time and memory usage of the AM model trained using FP32 and FP16 mixed precision
 3565 (FP16-mix), evaluated over 5 epochs with a training size of 10,000 in the CVRP20, CVRP50, and CVRP100.
 3566

Problem	Precision	Running time [s]	Memory usage [GiB]
CVRP20	FP32	6.33 ± 0.26	1.41 ± 0.04
	FP16-mix	5.89 ± 0.07	0.84 ± 0.01
CVRP50	FP32	13.58 ± 0.12	4.79 ± 0.40
	FP16-mix	11.68 ± 0.30	2.30 ± 0.25
CVRP100	FP32	35.09 ± 0.71	13.47 ± 0.63
	FP16-mix	25.11 ± 0.66	8.14 ± 0.82

3574
3575
3576 E.7.2 FLASHATTENTION
3577

3578 Given that the Attention operator is used on several occasions, especially in autoregressive models,
 3579 there is a need to support fast and efficient software routines that can compute this ubiquitous oper-
 3580 eration. In RL4CO, we natively support FlashAttention (Dao et al., 2022; Dao, 2023) from both
 3581 PyTorch 2.0+ and the original FlashAttention repository ²¹, to which we also made some minor
 3582 contributions when we found bugs.
 3583



3594
3595
3596
3597
3598 Figure 29: Running time of the graph attention encoder from the Attention Model, equipped with a standard
 3599 attention layer, FlashAttention1, and FlashAttention2, across different problem sizes for both (a) the TSP and
 3600 (b) the CVRP environments.
 3601

3602 As shown in Fig. 29, different implementations can make a difference, especially with large problem
 3603 sizes. It should be noted that while more scalable, FlashAttention at the moment is restricted to no
 3604 or causal masks only. Therefore, usage in the masked attention decoding scheme is not possible
 3605 at the moment, although it could be even more impactful due to the auto-regressive nature of our
 3606 encoder-decoder scheme. Recent works as Pagliardini et al. (2023) may be useful in extending
 3607 FlashAttention to other masking patterns. We note that masking should, in principle, be even faster
 3608 than un-masked attention, given that operations can be skipped in a per-block manner.
 3609

3610 E.7.3 EFFICIENT MEMORY HANDLING IN ENVIRONMENTS
 3611

3612 When dealing with RL problems, there is usually a tradeoff between memory and speed. This
 3613 happens because environments are parallelized using multiple processes or threads, the policy net-
 3614 work is replicated to each environment, or observations incoming from each environment need
 3615 to be gathered, sent to the policy network, and then the output action scattered back to the
 3616 representative environment. In the first case, network duplication causes large memory con-
 3617 sumption; in the second case, communication between processes slows down. In RL4CO, we
 solve the problem by using batched environments, i.e., every environment is responsible not

²¹Available at <https://github.com/Dao-AILab/flash-attention>.

for a single instance of a problem but a batch of instances at the same time. By doing so, the policy can live in the same process of the environment, in the same device, and receive and send batched data without any communication overhead or additional memory consumption. To further improve performances, we rewrite a core component of the TorchRL environment, namely the `step` method of the TorchRL base environment. The original `step` method performs some checks that, while useful for generic environments, can be omitted for RL4CO ones. It also duplicates the information in the output `TensorDict` by returning both the previous and the new state. In RL4CO, the previous state is always redundant, hence our `step` method does not keep it, reducing the memory consumption. We can see in Table 30 that using RL4CO step method has a great benefit in terms of speed, especially for high-dimensional environments. The results are collected for the TSP and CVRP environment during one epoch of training for a dataset of size 100000. The table shows the difference in training time and peak allocated memory for the training when the environment uses the TorchRL step method and the RL4CO step method. The peak allocated memory is computed using the `torch.cuda.max_memory_allocated` method from PyTorch, and experiments are run on a Tesla V100 DGX 32GB.

Table 30: Comparison of training time in seconds for one epoch with RL4CO and TorchRL step method.

Configuration		Step method	
Environment	Nodes	RL4CO	TorchRL
TSP	50	46.3	49.6
	100	102.9	108.6
	200	284.9	302.2
CVRP	50	72.9	73.4
	100	147.3	154.3
	200	371.7	406.4

E.8 TOWARDS FOUNDATION MODELS

Motivation Although learning to solve VRPs has gained significant attention, previous methods are only structured and trained independently on a specific problem, making them less generic and practical. Inspired by the recent success of foundation models in the language and vision domains, some works started to build foundation models for VRPs (Liu et al., 2024a; Zhou et al., 2024; Berto et al., 2024), aiming to solve a wide spectrum of problem variants using a single model. The main idea is to train a (large) model on diverse VRPs, which can be represented by a unified template. Typically, VRPs share several common attributes. For example, CVRP and VRPTW share the capacity attribute while only differing in the time window attribute. Therefore, a simple template could be a union set of attributes that exist in all VRP variants. By training on diverse VRP variants leveraging this unified representation, the foundation VRP model has the potential to efficiently and effectively solve any variant, making it a favorable choice versus traditional solvers (e.g., OR-Tools (Perron and Furnon, 2023)) in the future.

E.8.1 EXPERIMENTAL SETTING

For traditional solvers, we use HGS-PyVRP (Wouda et al., 2024), an open-source VRP solver based on the state-of-the-art HGS-CVRP (Vidal, 2022), and Google’s OR-Tools (Perron and Furnon, 2023), an open-source solver based on constraint programming for complex optimization problems, to solve all VRP variants considered in this study. Both baseline methods solve each instance on a single CPU core with a time limit of 10 and 20 seconds for instances with 50 and 100 nodes, respectively. We parallelize traditional solvers across 16 CPU cores as in (Kool et al., 2019a). For neural solvers, we mostly follow the training setups from previous works (Liu et al., 2024a; Zhou et al., 2024; Berto et al., 2024). In specific, the model is trained over 300 epochs, with each epoch containing 100,000 instances generated on the fly. The Adam optimizer is used with a learning rate of $3e - 4$, a weight decay of $1e - 6$, and a batch size of 256. The learning rate decays by 10 at 270 and 295 epochs. Note that different from Liu et al. (2024a); Zhou et al. (2024), we allow various problem variants to be trained in each batch training following Berto et al. (2024). We consider 16 VRP variants as shown in Table 8, including the constraints of capacity, time window, backhaul, open route, and duration limit. The training variants include CVRP, OVRP, VRPL, VRPB, VRPTW, and OVRPTW. During inference, we use greedy rollout with x8 instance augmentation following Kwon et al. (2020). We report the average results (i.e., objective values and gaps) over the test dataset that contains 1,000 instances, and the total time to solve the entire test dataset. The gaps are computed with respect to the results of HGS-PyVRP. All neural solvers are implemented using RL4CO.

Table 31: Performance on 1,000 test instances. * represents 0.000%, with which the gaps are computed.

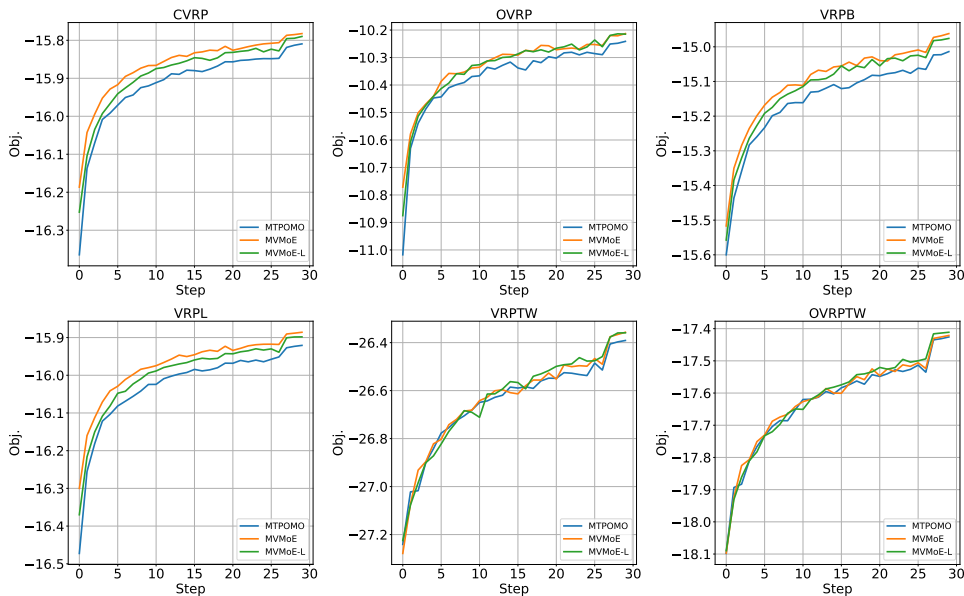
Method	N = 50			N = 100			Method	N = 50			N = 100			
	Obj.	Gap	Time	Obj.	Gap	Time		Obj.	Gap	Time	Obj.	Gap	Time	
CVRP	HGS-PyVRP	10.287	*	4.6m	15.543	*	9.2m	HGS-PyVRP	16.032	*	4.6m	25.433	*	9.2m
	OR-Tools	10.523	2.294%	4.6m	16.361	5.263%	9.2m	OR-Tools	16.124	0.574%	4.6m	25.923	1.927%	9.2m
	MTPOMO	10.408	1.176%	2s	15.809	1.711%	10s	MTPOMO	16.396	2.270%	2s	26.391	3.767%	11s
	MVMoE	10.397	1.069%	3s	15.782	1.538%	13s	MVMoE	16.394	2.258%	3s	26.357	3.633%	14s
	MVMoE-L	10.404	1.137%	3s	15.790	1.589%	12s	MVMoE-L	16.393	2.252%	3s	26.359	3.641%	13s
OVRP	HGS-PyVRP	6.494	*	4.6m	9.730	*	9.2m	HGS-PyVRP	10.328	*	4.6m	15.637	*	9.2m
	OR-Tools	6.555	0.939%	4.6m	10.081	3.607%	9.2m	OR-Tools	10.570	2.343%	4.6m	16.466	5.302%	9.2m
	MTPOMO	6.712	3.357%	2s	10.241	5.252%	10s	MTPOMO	10.454	1.220%	2s	15.921	1.816%	12s
	MVMoE	6.696	3.111%	3s	10.213	4.964%	13s	MVMoE	10.442	1.104%	3s	15.886	1.592%	13s
	MVMoE-L	6.704	3.234%	2s	10.215	4.985%	12s	MVMoE-L	10.450	1.181%	2s	15.898	1.669%	10s
VRPB	HGS-PyVRP	9.688	*	4.6m	14.386	*	9.2m	HGS-PyVRP	10.485	*	4.6m	16.900	*	9.2m
	OR-Tools	9.829	1.455%	4.6m	15.010	4.338%	9.2m	OR-Tools	10.497	0.114%	4.6m	17.023	0.728%	9.2m
	MTPOMO	9.975	2.962%	2s	15.014	4.365%	10s	MTPOMO	10.664	1.707%	2s	17.426	3.112%	11s
	MVMoE	9.954	2.746%	3s	14.962	4.004%	13s	MVMoE	10.665	1.717%	3s	17.421	3.083%	15s
	MVMoE-L	9.963	2.839%	2s	14.976	4.101%	11s	MVMoE-L	10.665	1.717%	2s	17.411	3.024%	14s
OVRPB	HGS-PyVRP	6.897	*	4.6m	10.304	*	9.2m	HGS-PyVRP	6.904	*	4.6m	10.310	*	9.2m
	OR-Tools	6.940	0.623%	4.6m	10.611	2.979%	9.2m	OR-Tools	6.949	0.652%	4.6m	10.613	2.939%	9.2m
	MTPOMO	7.392	7.177%	2s	11.787	14.392%	10s	MTPOMO	7.400	7.184%	2s	11.786	14.316%	10s
	MVMoE	7.566	9.700%	3s	11.873	15.227%	13s	MVMoE	7.577	9.748%	3s	11.875	15.179%	13s
	MVMoE-L	7.388	7.119%	2s	11.806	14.577%	12s	MVMoE-L	7.391	7.054%	2s	11.814	14.588%	12s
OVRPBLTW	HGS-PyVRP	11.597	*	4.6m	19.005	*	9.2m	HGS-PyVRP	11.590	*	4.6m	19.167	*	9.2m
	OR-Tools	11.612	0.129%	4.6m	19.198	1.016%	9.2m	OR-Tools	11.610	0.173%	4.6m	19.314	0.767%	9.2m
	MTPOMO	11.986	3.354%	2s	20.048	5.488%	11s	MTPOMO	11.980	3.365%	2s	20.209	5.436%	11s
	MVMoE	11.949	3.305%	3s	20.092	5.720%	15s	MVMoE	11.957	3.167%	3s	20.254	5.671%	15s
	MVMoE-L	11.961	3.139%	3s	20.033	5.409%	14s	MVMoE-L	11.951	3.115%	2s	20.173	5.249%	14s
OVRPL	HGS-PyVRP	6.510	*	4.6m	9.709	*	9.2m	HGS-PyVRP	10.455	*	4.6m	16.962	*	9.2m
	OR-Tools	6.571	0.937%	4.6m	10.047	3.481%	9.2m	OR-Tools	10.465	0.096%	4.6m	17.100	0.814%	9.2m
	MTPOMO	6.732	3.410%	2s	10.216	5.222%	10s	MTPOMO	10.625	1.626%	2s	17.486	3.089%	11s
	MVMoE	6.713	3.118%	3s	10.187	4.923%	13s	MVMoE	10.631	1.683%	3s	17.483	3.072%	15s
	MVMoE-L	6.725	3.303%	2s	10.185	4.903%	12s	MVMoE-L	10.635	1.722%	3s	17.474	3.019%	14s
VRPBL	HGS-PyVRP	9.688	*	4.6m	14.373	*	9.2m	HGS-PyVRP	18.361	*	4.6m	29.026	*	9.2m
	OR-Tools	9.820	1.363%	4.6m	15.084	4.947%	9.2m	OR-Tools	18.422	0.332%	4.6m	29.830	2.770%	9.2m
	MTPOMO	9.994	3.159%	2s	15.033	4.592%	10s	MTPOMO	19.028	3.633%	2s	31.062	7.014%	11s
	MVMoE	9.971	2.921%	3s	14.979	4.286%	13s	MVMoE	18.967	3.300%	3s	31.114	7.194%	15s
	MVMoE-L	9.977	2.983%	2s	14.990	4.293%	11s	MVMoE-L	18.998	3.469%	3s	31.032	6.911%	13s
VRPBLTW	HGS-PyVRP	18.167	*	4.6m	29.000	*	9.2m	HGS-PyVRP	15.951	*	4.6m	25.678	*	9.2m
	OR-Tools	18.374	1.139%	4.6m	29.964	3.324%	9.2m	OR-Tools	16.036	0.533%	4.6m	26.156	1.862%	9.2m
	MTPOMO	18.995	4.558%	2s	31.184	7.531%	11s	MTPOMO	16.310	2.251%	2s	26.650	3.785%	11s
	MVMoE	18.934	4.222%	3s	31.223	7.666%	15s	MVMoE	16.315	2.282%	3s	26.635	3.727%	14s
	MVMoE-L	18.970	4.420%	2s	31.138	7.372%	14s	MVMoE-L	16.311	2.257%	3s	26.637	3.735%	13s

E.8.2 EMPIRICAL RESULTS

We show the comprehensive evaluation results and validation curves in Table 31 and Fig. 30, respectively. The conclusions are consistent with previous studies (Liu et al., 2024a; Zhou et al., 2024; Berto et al., 2024) that 1) the foundation VRP solvers exhibit remarkable zero-shot generalization performance, even only trained on several VRPs with simple constraints; 2) conditional computation (e.g., mixture-of-experts (Jacobs et al., 1991; Shazeer et al., 2017)) can greatly enhance the model capacity without a proportional increase in computation. In Table 32, we further show the performance on CVRPLIB (Lima et al., 2014), which is a real-world benchmark dataset including instances with diverse distributions. We empirically observe that training on multiple VRPs can significantly improve the out-of-distribution generalization performance of neural VRP solvers, demonstrating the great promise of developing foundation models in VRPs.

E.8.3 DISCUSSION

Foundation models, a class of large-scale deep learning models pre-trained on extensive datasets of diverse tasks, have recently revolutionized the fields of language and vision domains. They can generate text, translate languages, summarize content, and more, all without task-specific training.

Figure 30: The validation curves of foundation models on $N = 100$.

This versatility makes them incredibly useful across various applications, from chatbots to academic research. Aiming for a more powerful and general solver, recent studies explore the possibility of pretraining a large model on a huge amount of optimization tasks. The long-term goal is to develop a foundation model for VRPs (or more broadly COPs), which can efficiently solve any problem variant, comparably or better to the conventional solvers with respect to the solution quality and inference speed. Despite the recent advancements of foundation VRP models (Liu et al., 2024a; Zhou et al., 2024; Berto et al., 2024), there are many challenges that need to be addressed by the NCO community, including but not limited to: 1) *scaling*: current autoregressive-based models are challenging to scale to the parameter levels of large language models (e.g., billions of parameters) due to the expensive training cost. RL-based training is data inefficient and converges slowly, whereas SL-based training requires a significant amount of optimal solutions, which are non-trivial to obtain for NP-hard problems. They also fail to be efficiently trained on large-scale instances; 2) *performance*: the empirical results are still far short of traditional solvers (e.g., OR-Tools). They may also suffer from generalization and robustness issues; 3) *generality*: the current problem formulation or template cannot solve novel problem variants in a zero-shot manner; 4) *interpretability*: the decision-making of foundation models is hard to explain.

Moreover, there is another line of research leveraging the existing large language models (LLMs) to generate solutions (Yang et al., 2024; Liu et al., 2023; Iklassov et al., 2024) or algorithms (Romera-Paredes et al., 2024; Liu et al., 2024b; Ye et al., 2024a), yielding impressive results when integrated with problem-specific heuristics or general meta-heuristics. Some studies employ LLMs to investigate the interpretability of solvers (Kikuta et al., 2024), automate problem formulation or simplify the use of domain-specific tools (Xiao et al., 2024; AhmadiTeshnizi et al., 2024; Wasserkrug et al., 2024) through text prompts. However, their performance is highly dependent on the utilized LLMs, and their outputs may be extremely sensitive to the designed prompts.

We view both as promising directions towards foundation models in combinatorial optimization. We call the attention from both the machine learning (ML) and operations research (OR) communities to advance the development of impactful foundational models and learning methods that are scalable, robust, generalizable, and interpretable across various optimization tasks in future work.

E.9 GENERALIZATION OF TRAINING ON MULTIPLE DISTRIBUTIONS AND MULTIPLE TASKS

Recent neural methods mostly train and test neural networks on the same task with instances of the same distribution and size, and hence suffer from inferior generalization performance. Some attempts have been made to alleviate the generalization issue, focusing on either distribution (Bi

3780 Table 32: Results on CVRPLib datasets with diverse distributions and sizes. All models are only trained on the
 3781 uniformly distributed data with the size $N = 100$.

3782

3783

Benchmark	Size N	Ins. Num.	POMO-CVRP		MTPOMO		MVMoE		MVMoE-L	
			Obj.	Gap	Obj.	Gap	Obj.	Gap	Obj.	Gap
3785 Set A	31-79	27	1088.5	4.9%	1084.2	4.3%	1081.0	3.8%	1085.4	4.4%
3786 Set B	30-77	23	1013.9	5.5%	1010.3	5.0%	1003.5	4.0%	1001.2	4.0%
3787 Set F	44-134	3	796.0	12.7%	812.7	16.3%	819.0	13.8%	799.0	14.1%
3788 Set M	100-199	5	1157.4	6.3%	1179.4	8.6%	1181.8	8.8%	1151.4	6.0%
3789 Set P	15-100	23	643.9	14.7%	621.8	8.4%	616.1	5.9%	619.8	6.9%
3790 Set X	100-1000	100	77199.6	21.1%	71153.8	11.7%	72798.7	15.0%	72446.1	13.9%

3791 [et al., 2022](#); [Jiang et al., 2022](#); [Xin et al., 2022](#)) or size ([Son et al., 2023](#)). More aligned to the
 3792 diverse distribution and size settings in the benchmark dataset TSPLib and CVRPLib, [Manchanda](#)
 3793 [et al. \(2023\)](#) and [Zhou et al. \(2023\)](#) consider generalization across both distribution and size in VRPs.

3794 However, these generalization methods adopt extra model architectures and training paradigms, re-
 3795 sulting in additional computational burdens. As a more efficient alternative, we observe that diversi-
 3796 fied training datasets significantly improve generalization performance. Specifically, as indicated in
 3797 the prior works, training on mixed distributions ([Bi et al., 2022](#)) and mixed VRP variants ([Liu et al.,](#)
 3798 [2024a](#); [Zhou et al., 2024](#); [Berto et al., 2024](#)) boosts the generalization capability. RL4CO, detailed
 3799 in [Appendix B.1.6](#), supports multiple VRP variants and the generation of diverse coordinate distri-
 3800 butions, enabling straightforward experimental setups. The implementation specifics are outlined in
 3801 [Appendix D.3.4](#). Evaluation results on the CVRPLib [Lima et al. \(2014\)](#), summarized in [Table 5](#) and
 3802 fully detailed in [Table 33](#), demonstrate that training across multiple distributions (i.e., MDPOMO)
 3803 achieves better generalization on datasets of similar size to the training set, whereas training across
 3804 multiple VRP tasks (i.e., MTPOMO) exhibits superior generalization across larger and more diverse
 3805 distributions. This indicates that different VRP variants share foundational knowledge, and learning
 3806 from this diversity enhances generalization beyond conventional training on a single distribution,
 3807 size, and task. These key findings highlight the necessity of developing foundational models across
 3808 diverse combinatorial optimization domains.

3809
3810
3811
3812
3813
3814
3815
3816
3817
3818
3819
3820
3821
3822
3823
3824
3825
3826
3827
3828
3829
3830
3831
3832
3833

