





# SVRPBench: A Realistic Benchmark for Stochastic Vehicle Routing Problem

Ahmed Heakl<sup>1\*</sup> Yahia Salaheldin Shaaban<sup>1\*</sup>  
Martin Takáč<sup>1</sup> Salem Lahlou<sup>1</sup> Zangir Iklassov<sup>1</sup>

<sup>1</sup>MBZUAI, Abu Dhabi, UAE

 <https://github.com/yehias21/vrp-benchmarks>  
 <https://huggingface.co/datasets/MBZUAI/svrp-bench>

## Abstract

Robust routing under uncertainty is central to real-world logistics, yet most benchmarks assume static, idealized settings. We present SVRPBench, the first open benchmark to capture high-fidelity stochastic dynamics in vehicle routing at urban scale. Spanning more than 500 instances with up to 1000 customers, it simulates realistic delivery conditions: time-dependent congestion, log-normal delays, probabilistic accidents, and empirically grounded time windows for residential and commercial clients. Our pipeline generates diverse, constraint-rich scenarios, including multi-depot and multi-vehicle setups. Benchmarking reveals that state-of-the-art RL solvers like POMO and AM degrade by over 20% under distributional shift, while classical and metaheuristic methods remain robust. To enable reproducible research, we release the dataset ([Hugging Face](#)) and evaluation suite ([GitHub](#)). SVRPBench challenges the community to design solvers that generalize beyond synthetic assumptions and adapt to real-world uncertainty.

## 1 Introduction

Efficient vehicle routing is fundamental to modern logistics and last-mile delivery. The classical Vehicle Routing Problem (VRP) [8, 11] seeks cost-effective routes for servicing customers under constraints such as vehicle capacities and time windows. Although well studied, real-world deployments face uncertain and dynamic conditions that most existing benchmarks do not adequately capture.

One key extension addressing real-world complexity is the *Stochastic Vehicle Routing Problem* (SVRP). Unlike deterministic VRP, SVRP explicitly incorporates uncertainty into routing decisions, with problem elements such as travel times, customer demands, service times, and even customer presence considered random variables [11, 22]. Consequently, routes are planned *a priori*, and corrective actions, known as recourse strategies, are applied when realized conditions deviate from planned values [9, 2]. Prominent examples include random travel times modeled by probabilistic distributions or random customer presence known as probabilistic VRP (PVRP) [18, 5]. Despite this extensive body of research, many existing public benchmarks for SVRP still rely on static assumptions, such as deterministic travel times, fixed customer availability, and unchanged route constraints, thus limiting their practical applicability and robustness evaluations, as shown in Table 1.

---

\*Equal contribution.

Table 1: Comparison of SVRPBench with existing VRP benchmarks. ✓ indicates full support, △ indicates partial or limited support, and ✗ indicates no support.

Feature	SVRPBench	CVRPLIB	SINTEF	VRP-REP	TSPLIB	RL4CO
<i>Stochastic Elements</i>						
Time-dependent travel delays	✓	✗	△	△	✗	✗
Peak-hour traffic patterns	✓	✗	✗	✗	✗	✗
Random travel time noise	✓	✗	△	△	✗	△
Probabilistic accidents	✓	✗	✗	✗	✗	✗
Heterogeneous time windows	✓	✗	△	△	✗	✗
<i>Problem Configurations</i>						
Multi-depot support	✓	△	✓	✓	✗	✗
Multi-vehicle fleets	✓	✓	✓	✓	✗	✓
Capacity constraints	✓	✓	✓	✓	✗	✓
Time window constraints	✓	△	✓	✓	✗	△
Clustered customer distributions	✓	△	△	✓	△	✗
<i>Scale &amp; Diversity</i>						
Small instances ( $\leq 100$ customers)	✓	✓	✓	✓	✓	✓
Medium instances (100-300)	✓	✓	✓	✓	△	✓
Large instances ( $>300$ )	✓	△	△	△	✗	△
Varying stochasticity levels	✓	✗	△	△	✗	✗

**The Case for a Realistic SVRP Benchmark.** Urban logistics operates under dynamic and uncertain conditions, yet most existing benchmarks fail to reflect this complexity. Practical routing systems must account for peak-hour congestion, random incidents like accidents, and diverse delivery preferences across customer types [14, 3, 24]. Ignoring these factors leads to overly optimistic performance assessments and misdirects algorithmic development toward unrealistic assumptions [1].

**Our Contributions.** To address these gaps, we introduce SVRPBench, a novel, open-source benchmark suite for the Stochastic Vehicle Routing Problem (SVRP), designed to simulate realistic logistics scenarios with embedded uncertainty. Our key contributions include:

- **Stochastic Realism.** We model time-dependent congestion using Gaussian mixtures, inject log-normal delays and probabilistic accidents [18], and generate customer time windows from empirical residential and commercial distributions.
- **Constraint-Rich Instance Generation.** Our framework supports multi-depot and multi-vehicle setups, strict capacity constraints, and diverse time window widths, all grounded in spatially realistic demand distributions.
- **Diverse Baseline Evaluation.** We benchmark classical heuristics (e.g., Nearest Neighbor, 2-opt), metaheuristics (e.g., ACO, Tabu Search [12, 7]), industrial solvers (OR-Tools [26], LKH3 [31]), and learning-based methods (AM [15], POMO [17]), highlighting how stochastic conditions affect solution quality, feasibility, and robustness.
- **Open Community Platform.** We release datasets, solvers, and evaluation scripts through a public repository to support reproducibility and foster future contributions.

By advancing realism and accessibility in SVRP benchmarking, SVRPBench aims to accelerate the development of robust, deployable routing algorithms suited for real-world logistics.

## 2 Realistic Stochastic Modeling

A core contribution of SVRPBench is its simulation of real-world uncertainty in urban-scale logistics. Classical VRP benchmarks often assume static travel times and rigid customer schedules [13], overlooking time-varying conditions and operational stochasticity. Informed by empirical and theoretical literature [3, 14, 1, 23, 25, 27, 19, 6, 10, 20], our benchmark introduces: (1) time-dependent congestion, (2) stochastic travel time delays, (3) accident-induced disruptions, and (4) customer-specific time window distributions.

## 2.1 Time-Dependent Travel Time Modeling

We model the travel time from node  $a$  to  $b$  at time  $t$  as:

$$T(a, b, t) = \frac{D(a, b)}{V} + B(a, b, t) \cdot R(t) + I_{\text{accidents}}(t) \cdot D_{\text{accident}}, \quad (1)$$

where  $D(a, b)$  is Euclidean distance and  $V$  is average road speed. The congestion factor  $B(a, b, t)$  is defined as:

$$B(a, b, t) = \alpha \cdot F_{\text{time}}(t) \cdot F_{\text{distance}}(D(a, b)), \quad (2)$$

with:

$$F_{\text{time}}(t) = \beta + \gamma \cdot [f(t; \mu_{\text{morning}}, \sigma_{\text{peak}}) + f(t; \mu_{\text{evening}}, \sigma_{\text{peak}})], \quad (3)$$

$$f(t; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{t-\mu}{\sigma}\right)^2}, \quad (4)$$

$$F_{\text{distance}}(D) = 1 - e^{-D/\lambda_{\text{dist}}}, \quad (5)$$

where the Gaussian peaks around  $\mu_{\text{morning}} = 8$  and  $\mu_{\text{evening}} = 17$  ( $\sigma_{\text{peak}} = 1.5$ ) align with observed urban traffic congestion patterns [27]. The distance decay  $\lambda_{\text{dist}} = 50$  modulates slowdown severity, reflecting empirical findings that longer trips are more likely to encounter congestion [6].

The multiplicative stochastic delay  $R(t)$  is drawn from a log-normal distribution:

$$\mu(t) = \mu_{\text{base}} + \delta \cdot [f(t; \mu_{\text{morning}}, \sigma_{\text{peak}}) + f(t; \mu_{\text{evening}}, \sigma_{\text{peak}})], \quad (6)$$

$$\sigma(t) = \sigma_{\text{base}} + \epsilon \cdot [f(t; \mu_{\text{morning}}, \sigma_{\text{peak}}) + f(t; \mu_{\text{evening}}, \sigma_{\text{peak}})], \quad (7)$$

$$R(t) \sim \text{LogNormal}(\mu(t), \sigma(t)), \quad (8)$$

reflecting both the skewed and bursty nature of traffic delays [19, 6]. Baseline values  $\mu_{\text{base}} = 0$  and  $\sigma_{\text{base}} = 0.3$  reflect free-flow conditions, while  $\delta = 0.1$  and  $\epsilon = 0.2$  capture peak-hour amplification.

Accident delays are modeled using a time-inhomogeneous Poisson process:

$$\lambda(t) = \lambda_{\text{scale}} \cdot f(t; \mu_{\text{night}}, \sigma_{\text{acc}}), \quad (9)$$

$$I_{\text{accidents}}(t) \sim \text{Poisson}(\lambda(t)), \quad (10)$$

$$D_{\text{accident}} \sim U(d_{\text{min}}, d_{\text{max}}), \quad (11)$$

where accidents peak around  $\mu_{\text{night}} = 21$  ( $\sigma_{\text{acc}} = 2$ ) due to elevated nighttime risks from fatigue and impaired driving [28]. The delay duration is drawn from  $U(0.5, 2.0)$  hours, consistent with industry reports on incident clearance times [28].

## 2.2 Customer Time Window Sampling

Residential and commercial customers exhibit different temporal availability patterns [23, 20]. For residential profiles, delivery windows are sampled from a bimodal Gaussian mixture:

$$T_{\text{start}} \sim \begin{cases} \mathcal{N}(\mu_{\text{res,morning}}, \sigma_{\text{res,morning}}^2), & \text{w.p. } 0.5, \\ \mathcal{N}(\mu_{\text{res,evening}}, \sigma_{\text{res,evening}}^2), & \text{w.p. } 0.5, \end{cases} \quad (12)$$

where  $\mu_{\text{res,morning}} = 480$  (8:00 AM) and  $\mu_{\text{res,evening}} = 1140$  (7:00 PM), with variances  $\sigma = 90$  and 120 mins, respectively, aligning with common parcel service offerings such as FedEx and Bring [10, 20]. The window duration is drawn from:

$$W_{\text{length}} \sim U(w_{\text{min}}, w_{\text{max}}), \quad T_{\text{start}} = \max(0, \min(T_{\text{start}}, 1440 - W_{\text{length}})). \quad (13)$$

Commercial customers follow a single-mode Gaussian:

$$T_{\text{start}} \sim \mathcal{N}(\mu_{\text{com}}, \sigma_{\text{com}}^2), \quad W_{\text{length}} \sim U(w_{\text{min}}, w_{\text{max}}^{\text{com}}), \quad (14)$$

with  $\mu_{\text{com}} = 780$  (1:00 PM),  $\sigma_{\text{com}} = 60$ , and  $w_{\text{max}}^{\text{com}} = 120$  minutes, reflecting standard daytime business hours and delivery norms [29].

This probabilistic windowing model encourages algorithms to balance varied service constraints, simulating realistic scheduling trade-offs in last-mile delivery systems.

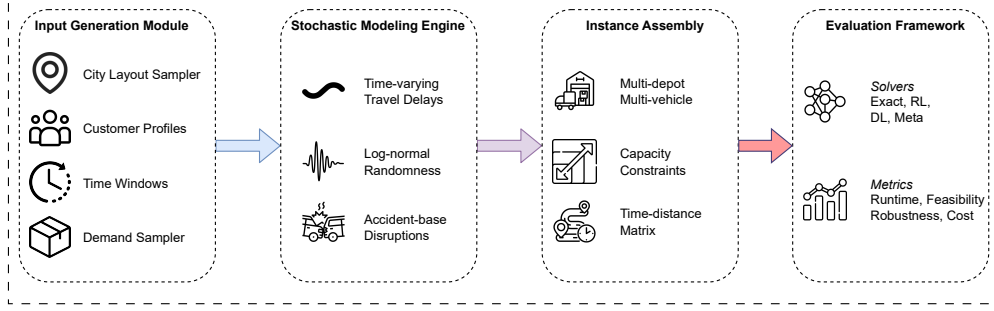


Figure 1: **SVRPBench pipeline.** The framework generates realistic SVRP instances through four stages: input generation, stochastic modeling, instance assembly, and evaluation with standardized metrics and solvers.

### 3 Dataset Construction Pipeline

To enable scalable and reproducible experimentation, we develop a unified pipeline that generates diverse, constraint-rich SVRP instances grounded in stochastic realism. It integrates models of customer behavior, traffic patterns, spatial layouts, and routing constraints to produce problem scenarios suited for evaluating both classical and learning-based solvers under realistic uncertainty [22, 11]. The complete pipeline is illustrated in Figure 1.

**Location Sampling.** We begin by selecting the total number of customers from  $\{10, 20, 100, 500, 1000\}$ , then compute the number of cities as  $\max(1, \#customers/50)$ . To simulate spatial separation between urban clusters, we apply K-Means clustering to generate city centers that are as distant from each other as possible. Customer locations are then sampled around each city center using 2D Gaussian distributions [14].

**Demand Assignment.** Each customer is assigned a discrete demand selected uniformly at random from a set  $\{1, 2, \dots, \max\_demand\}$ . The number of vehicles and their capacity are computed based on the total customer demand, with vehicle capacity set as  $\text{total demand} \div \text{number of vehicles}$ . This ensures balanced feasibility across instance scales [9].

**Time Window Assignment.** Customer time windows are generated stochastically, following the models described in Section 2. Residential and commercial customer patterns are differentiated using realistic temporal distributions [3].

**Travel Time Matrix Construction.** A full travel time matrix  $T(a, b, t)$  is computed for all location pairs, incorporating deterministic base time, time-dependent congestion patterns, log-normal stochastic variation, and random accident delays, as detailed in Section 2. This captures the nonlinear, time-varying nature of urban transportation systems [18].

**Constraint Integration.** We support both single-depot and multi-depot configurations. In multi-depot settings, depots can be placed either randomly or aligned with city centers (one per city). A homogeneous fleet of vehicles is used, and vehicle count is configured to balance demand and capacity. All customer time windows are sampled to ensure feasibility under the assigned travel time model [1].

**Validation.** Each generated instance undergoes automated validation to ensure feasibility under both capacity and temporal constraints. For CVRP, we verify that the total vehicle capacity (number of vehicles  $\times$  per-vehicle capacity) exceeds the sum of all customer demands, ensuring that a feasible route covering all customers exists. For TWVRP, we construct a time-windowed demand histogram by binning the time axis and accumulating customer demands per bin. We then identify the peak-demand bin and ensure that the fleet capacity is sufficient to serve this worst-case demand, i.e.,  $\text{capacity} \times \text{num\_vehicles} \geq \max_t \text{demand}(t)$ . This provides a conservative guarantee that even

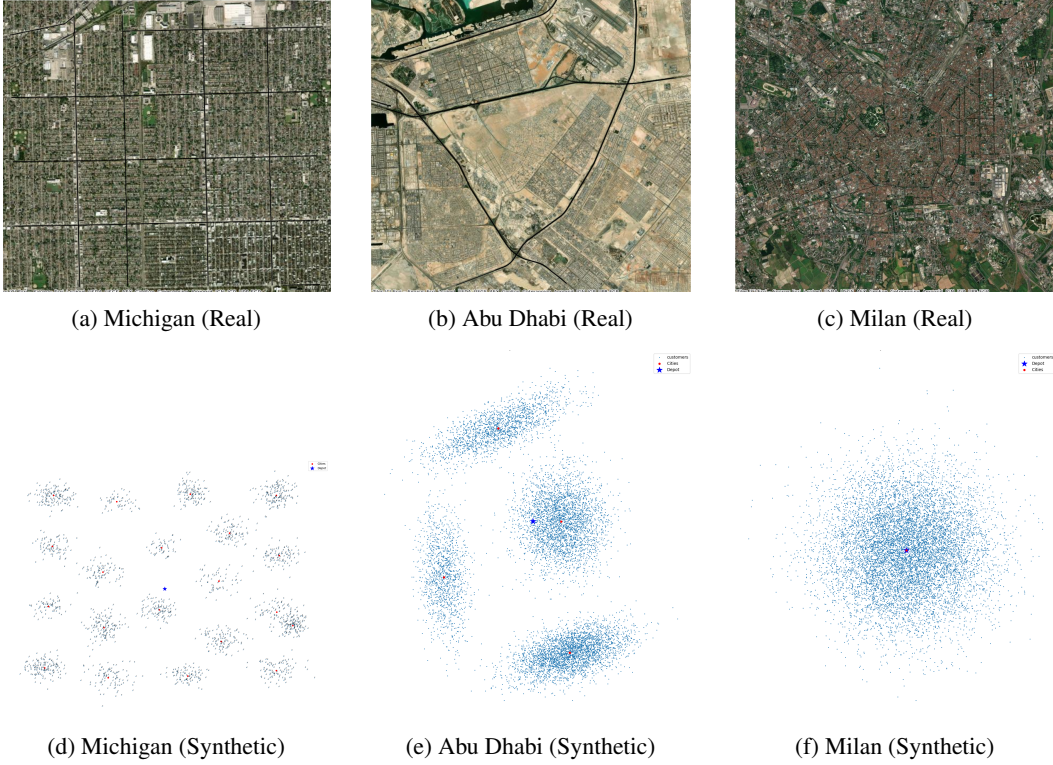


Figure 2: Comparison of real (top) and synthetic (bottom) routing instances across three cities.

under concentrated temporal demand, a feasible schedule remains possible. Infeasible instances (e.g., unreachable nodes or incompatible time windows) are filtered or regenerated.

Parameters are selected to reflect urban-scale routing challenges but can be modified for rural or industrial scenarios. Accident frequency and delay magnitudes are parameterized using a Poisson-based arrival model and uniform delay range, respectively. Customer types are split roughly 60% residential to 40% commercial, matching empirical logistics patterns [3].

**Various Scales.** Our benchmark includes three instance tiers. *Small* instances (50–100 customers, 1–2 depots) with low noise allow quick testing. *Medium* instances (100–300 customers, 2–3 depots) feature moderate stochasticity. *Large* instances (300+ customers) integrate high travel-time variability and tighter delivery windows to stress-test scalability. All levels are generated with multiple random seeds to support statistical averaging and ensure robustness of comparisons.

To validate the realism of our spatial sampling strategy, we visually compare synthetic routing instances against satellite imagery of real-world cities. As shown in Figure 2, our generated layouts closely mimic key structural patterns, grid-like in Michigan, radial in Milan, and dispersed in Abu Dhabi, demonstrating the pipeline’s ability to emulate diverse urban morphologies critical for evaluating routing algorithms in geographically grounded scenarios.

## 4 Evaluation Protocol

To ensure fair, rigorous, and reproducible comparisons across routing algorithms, we propose a standardized evaluation protocol tailored for our stochastic vehicle routing benchmark. This protocol assesses not only solution quality but also robustness, feasibility, and scalability under conditions of realistic uncertainty, addressing limitations of earlier benchmark designs that overlooked stochastic effects [22, 11].

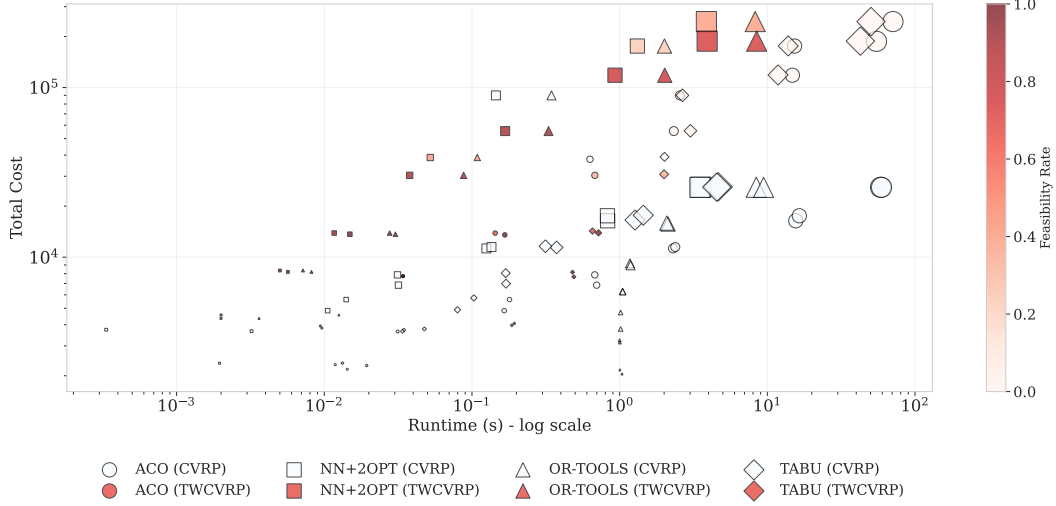


Figure 3: Solver Comparison: Overall Performance Metrics.

#### 4.1 Performance Metrics

We report a comprehensive suite of metrics to evaluate different facets of algorithmic behavior. The *Total Cost (TC)* measures the cumulative travel time across all vehicles, including congestion-induced delays and accident-based disruptions. Formally, it is computed as:

$$TC = \sum_{k \in V} \sum_{(i,j) \in \text{route}_k} T(i, j, t_i), \quad (15)$$

where  $T(i, j, t_i)$  is the sampled travel time from node  $i$  to  $j$  at time  $t_i$ .

*Constraint Violation Rate (CVR)* quantifies the proportion of customers whose service violates time windows or exceeds vehicle capacity, capturing solution feasibility:

$$CVR = \frac{\# \text{violations}}{\# \text{customers}} \times 100\%. \quad (16)$$

*Feasibility Rate (FR)* reflects the robustness of solutions across instances and solvers. It is defined as the fraction of problem instances for which a solution satisfies all routing constraints:

$$FR = \frac{\# \text{feasible instances}}{\# \text{total instances}}. \quad (17)$$

*Runtime (RT)* captures wall-clock computation time, serving as a proxy for scalability and practical deployability.

*Robustness (ROB)* measures the variability in cost due to stochastic elements by computing the variance across  $N$  independent samples of the same instance:

$$ROB = \frac{1}{N} \sum_{i=1}^N (TC_i - \overline{TC})^2, \quad (18)$$

where  $\overline{TC}$  denotes the mean total cost. This metric is especially important in stochastic VRP settings [2, 18].

## 5 Experimental Results

We conduct a comprehensive evaluation of baseline methods on our stochastic VRP benchmark, which systematically varies four key dimensions: instance size, problem type, depot configuration, and vehicle configuration.

We generate 10 instances for each combination across instance sizes {10, 20, 50, 100, 200, 500, 1000}, problem types {CVRP, TWVRP}, depot configurations {single, multi}, and vehicle settings

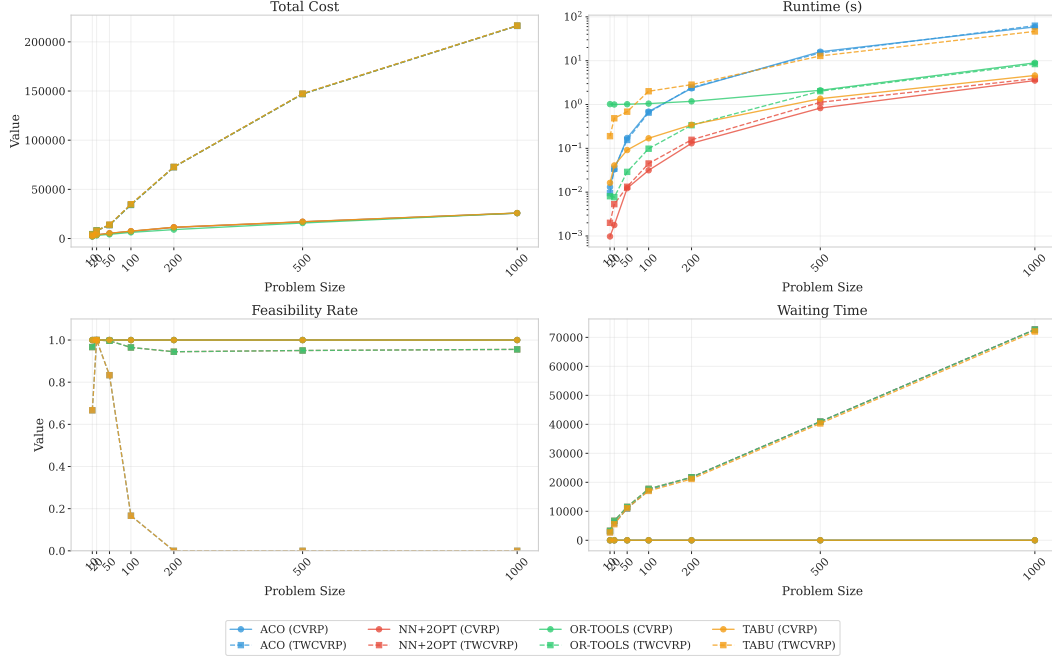


Figure 4: Solver Performance by Problem Size.

{single, multi}, yielding a large-scale, structured test suite. Additionally, we provide a scalable data generator for training. Reinforcement learning models were trained on 100k synthetic instances under the single-depot, single-vehicle CVRP and TWCVRP regimes.

## 5.1 Evaluation Scope

All methods were evaluated under the stochastic setting defined in Section 2. Metrics reported include total cost (incorporating all stochastic factors), constraint violation rate (CVR), feasibility rate, runtime, and robustness (measured as variance across stochastic samples).

Classical algorithms, Nearest Neighbor + 2-opt, Tabu Search, and ACO (refer to Appendix B for more details), were evaluated across all settings without modification. Their flexibility allows them to handle diverse configurations out of the box.

## 5.2 Experimental Setup

All baselines were evaluated on a consumer-grade CPU (Intel i7, 16GB RAM), except learning-based models, which used a single NVIDIA RTX 4080. Classical and metaheuristic solvers were implemented in Python; learning models used the RL4CO framework [4]. Training for RL models was done on 100k synthetic instances (refer to Appendix D for more details). Evaluation followed the stochastic protocol detailed in Section 2, averaging results over five realizations per test case.

## 5.3 Results & Analysis

**Overall Performance.** Table 2 and Figure 3 summarize the aggregate performance across all test cases. OR-Tools achieved the best overall cost (40,259), followed closely by ACO (40,566; +0.8%) and POMO (40,650; +1.0%), with OR-Tools and NN+2opt maintaining the highest feasibility rates (98.4%) while NN+2opt delivered the fastest runtime (0.697s). Learning-based approaches demonstrated a feasibility-speed tradeoff, with POMO offering better solution quality than NN+2opt at competitive runtimes (1.421s) while the Attention Model showed higher constraint violations (CVR: 1.9%) but reasonable performance across other metrics.



Table 2: Performance of baseline methods (mean over all instances, 5 stochastic runs).

Method	Total Cost ↓	CVR (%) ↓	Feasibility ↑	Runtime (s) ↓	Robustness ↓
NN+2opt	40707.5	1.6	0.984	0.697	0.1
Tabu Search	40787.8	1.6	0.690	5.157	0.1
ACO	40566.5	1.6	0.690	11.382	0.1
OR-Tools	40259.3	1.6	0.984	1.940	0.1
Attention Model (AM)	41358.3	1.9	0.910	1.852	0.2
POMO	40650.4	1.7	0.933	1.421	0.1

Table 3: Performance Comparison: CVRP vs TWCVRP.

Method	CVRP				TWCVRP				Impact % Δ
	Cost↓	CVR↓	Feas↑	RT↓	Cost↓	CVR↓	Feas↑	RT↓	
NN+2opt	10399.2	0.0	1.000	646.3	71015.8	3.2	0.968	747.8	+582.9
Tabu Search	10494.1	0.0	1.000	945.1	71081.5	3.2	0.381	9368.6	+577.3
ACO	10384.9	0.0	1.000	11159.8	70748.1	3.2	0.381	11603.6	+581.3
OR-Tools	9499.7	0.0	1.000	2328.0	71018.8	3.2	0.968	1552.1	+647.6
Attention Model (AM)	11235.6	0.2	0.965	1775.4	71481.0	3.6	0.854	1929.2	+536.2
POMO	10358.7	0.1	0.987	1316.9	70942.1	3.3	0.879	1525.3	+584.8

**Impact of Time Windows.** Table 3 reveals that introducing time windows (TWCVRP) increases total cost by 536–648% across all solvers, with OR-Tools incurring the highest relative penalty (+647.6%) while the Attention Model showed the lowest relative increase (+536.2%). Learning-based methods demonstrated moderate resilience to time constraints with POMO maintaining 87.9% feasibility and Attention Model 85.4%, positioning them between the top performers (NN+2opt and OR-Tools at >96%) and the struggling metaheuristics (ACO and Tabu Search at 38.1%).

**Scalability by Instance Size.** As shown in Table 4 and Figure 4, cost scaled approximately 16× from small ( $\leq 50$  nodes) to large ( $\geq 500$  nodes) instances across all methods, with NN+2opt and OR-Tools maintaining feasibility >97% at all scales, while learning-based methods showed moderate degradation (POMO: 86%, AM: 83.5%). Learning-based approaches demonstrated competitive performance-runtime tradeoffs, with POMO offering the fastest runtime on small instances (29.7s) and maintaining feasibility significantly better than ACO and Tabu Search (50%) on large instances, though traditional heuristics still held the advantage for the largest problems.

**Effect of Depot Configuration.** Table 5 shows that multi-depot setups consistently reduced costs and improved feasibility across all methods, with OR-Tools achieving a 72% cost reduction (from 34,611 to 9,561) and POMO showing similarly impressive gains (71% reduction to 10,178). Learning-based methods particularly benefited from multi-depot configurations, with both POMO and Attention Model reaching perfect feasibility (100%) despite their variable performance in single-depot scenarios (92-96.5%), supporting the counterintuitive finding that more flexible depot placements improve both computational and solution efficiency regardless of algorithm class.

**Key Takeaways.** Our evaluation underscores several important insights:

- OR-Tools is the most reliable choice for large-scale offline optimization, balancing quality and feasibility despite higher runtimes.

Table 4: Detailed Performance Analysis by Instance Size.

Method	Small ( $\leq 50$ )				Medium (100-200)				Large ( $\geq 500$ )			
	Cost↓	CVR↓	Feas↑	RT↓	Cost↓	CVR↓	Feas↑	RT↓	Cost↓	CVR↓	Feas↑	RT↓
NN+2opt	6295.0	0.6	0.994	5.9	31486.1	2.3	0.977	90.9	101547.5	2.4	0.976	2340.0
Tabu Search	6232.5	0.6	0.917	251.6	31692.2	2.3	0.542	1339.5	101716.5	2.4	0.500	16332.1
ACO	6080.7	0.6	0.917	69.6	31371.9	2.3	0.542	1530.6	101490.0	2.4	0.500	38201.0
OR-Tools	6008.1	0.6	0.994	513.7	30640.2	2.3	0.977	665.8	101255.0	2.4	0.976	5353.7
Attention Model (AM)	6523.2	0.8	0.975	42.3	32165.5	2.6	0.910	857.4	102756.2	2.9	0.835	4758.9
POMO	6176.4	0.7	0.985	29.7	31024.8	2.4	0.945	642.3	101408.7	2.5	0.860	3586.2



Table 5: Performance Analysis by Depot Configuration.

Method	Single Depot				Multi Depot			
	Cost ↓	CVR ↓	Feas ↑	RT ↓	Cost ↓	CVR ↓	Feas ↑	RT ↓
NN+2opt	34978.5	0.8	0.992	686.3	10625.2	0.0	1.000	643.7
Tabu Search	35072.0	0.8	0.690	4818.2	10713.8	0.0	1.000	946.1
ACO	34852.1	0.8	0.690	10712.0	10614.9	0.0	1.000	11298.7
OR-Tools	34611.0	0.8	0.992	1911.2	9561.4	0.0	1.000	2396.5
Attention Model (AM)	35825.6	1.1	0.920	1785.3	10974.7	0.0	1.000	1852.6
POMO	34786.3	0.9	0.965	1438.2	10178.5	0.0	1.000	1324.8

- NN+2opt offers a robust, low-latency alternative for real-time deployment with minimal compromise on cost or feasibility.
- Metaheuristics underperform at scale, while learning-based methods like POMO offer feasible solutions with better scalability, though still lag behind top heuristics.
- The Attention Model demonstrates potential but requires further refinement to match the performance of top-performing methods, particularly for large instances.
- Time windows impose the most significant complexity, sharply degrading performance for non-adaptive solvers, though learning-based methods show moderate resilience.
- Multi-depot settings improve both feasibility and runtime across all solver types, offering a practical design consideration for logistics planning.

Together, Figures 3 and 4 illustrate these trends across key metrics. SVRPBench successfully reveals scalability bottlenecks, constraint sensitivity, and performance trade-offs, establishing a realistic and informative testbed for stochastic routing research. Please refer to appendix C for additional results.

## 6 Limitations and Future Directions

While SVRPBench advances realism in stochastic vehicle routing, several limitations remain. Our delay models rely on Gaussian and log-normal distributions to simulate traffic peaks and randomness—efficient and interpretable, yet unable to capture network-level dynamics such as bottlenecks, cascading congestion, or real-time rerouting [14]. These assumptions, however, are user-modifiable, allowing injection of domain-specific uncertainty. Reinforcement learning methods like AM and POMO show limited scalability to larger instances, reflecting overfitting and weak generalization. Additionally, our current evaluation protocol lacks standardized procedures to assess robustness across instance scales and distribution shifts, motivating future work on curriculum learning and hierarchical solver design.

To further bridge the gap to real-world logistics, future extensions will incorporate road-constrained instances derived from OpenStreetMap or GIS data, enabling geographically grounded routing behavior. Dynamic and multi-day settings—with online updates and rolling horizons—will support evaluation of adaptive strategies [2]. We also plan to introduce diagnostic tasks for probing model robustness, generalization under distributional shift, and few-shot performance [21, 17], enabling more fine-grained analysis of algorithmic reliability in complex environments.

## 7 Conclusion

We presented SVRPBench, a modular and open-source benchmark for evaluating vehicle routing under realistic stochastic dynamics. By incorporating time-dependent congestion, probabilistic delays, and heterogeneous customer time windows, our benchmark departs from static assumptions and reflects the operational uncertainty of real logistics.

Empirical results across over 500 instances revealed that classical and metaheuristic methods remain competitive on feasibility and runtime, while reinforcement learning models like POMO and AM, despite strong performance in training regimes, struggled with multi-depot generalization and exhibited  $> 20\%$  cost degradation under distributional shift. Surprisingly, multi-depot configurations consistently improved both cost and robustness, even for learning-based solvers, highlighting the importance of flexible depot placement in practical settings.

By supporting large-scale, reproducible evaluations via Hugging Face and GitHub, SVRPBench offers a community platform to benchmark solvers across realism axes. We urge the research community to develop adaptive, noise-aware routing algorithms that bridge the gap between synthetic optimization and deployable, resilient logistics solutions.

## References

- [1] Yossiri Adulyasak and Patrick Jaillet. Models and algorithms for stochastic and robust vehicle routing with deadlines. *Transportation Science*, 50(2):608–626, 2016.
- [2] Cock Bastian and Alexander H. G. Rinnooy Kan. The stochastic vehicle routing problem revisited. *European Journal of Operational Research*, 56(3):407–412, 1992.
- [3] Russell Bent and Pascal Van Hentenryck. Scenario-based planning for partially dynamic vehicle routing with stochastic customers. *Operations Research*, 52(6):977–987, 2004.
- [4] F. Berto, C. Hua, J. Park, M. Kim, H. Kim, J. Son, H. Kim, J. Kim, and J. Park. RL4co: A unified reinforcement learning for combinatorial optimization library. In *Proceedings of Advances of Neural Information Processing Systems (workshop)*, 2023.
- [5] Dimitris J. Bertsimas, Patrick Jaillet, and Amedeo R. Odoni. A priori optimization. *Operations Research*, 38(6):1019–1033, 1990.
- [6] Werner Brilon, Jürgen Geistefeldt, and Markus Regler. Reliability of travel times: A stochastic modeling approach. *Transportation Research Record*, 2061(1):1–8, 2008.
- [7] K. Chepuri and T. Homem-de Mello. Solving the vehicle routing problem with stochastic demands using the cross-entropy method. *Annals of Operations Research*, 134(1):153–181, 2005.
- [8] George B Dantzig and John H Ramser. The truck dispatching problem. *Management science*, 6(1):80–91, 1959.
- [9] Moshe Dror, Gilbert Laporte, and Pierre Trudeau. Vehicle routing with stochastic demands: Properties and solution frameworks. *Transportation Science*, 23(3):166–176, 1989.
- [10] FedEx Corporation. Fedex residential delivery options whitepaper. *Whitepaper*, 2020. Flexible delivery time window practices.
- [11] Michel Gendreau, Gilbert Laporte, and Renaud Séguin. Stochastic vehicle routing. *European Journal of Operational Research*, 88(1):3–12, 1996.
- [12] Michel Gendreau, Gilbert Laporte, and Renaud Séguin. A tabu search heuristic for the vehicle routing problem with stochastic demands and customers. *Operations Research*, 44(3):469–477, 1996.
- [13] Michel Gendreau, Gilbert Laporte, and Rene Seguin. Stochastic vehicle routing. *European Journal of Operational Research*, 88(1):3–12, 1996.
- [14] Lars Magnus Hvattum, Arne Lø kketangen, and Gilbert Laporte. Solving a dynamic and stochastic vehicle routing problem with a sample scenario hedging heuristic. *Transportation Science*, 40(4):421–438, 2006.
- [15] Wouter Kool, Herke Van Hoof, and Max Welling. Attention, learn to solve routing problems! *arXiv preprint arXiv:1803.08475*, 2018.
- [16] Wouter Kool, Herke van Hoof, and Max Welling. Attention, learn to solve routing problems! *International Conference on Learning Representations (ICLR)*, 2019.
- [17] Yeong-Dae Kwon, Jinho Choo, Byoungjip Kim, Iljoo Yoon, Youngjune Gwon, and Seungjai Min. Pomo: Policy optimization with multiple optima for reinforcement learning. *Advances in Neural Information Processing Systems*, 33:21188–21198, 2020.
- [18] Gilbert Laporte, François V. Louveaux, and Hélène Mercure. The vehicle routing problem with stochastic travel times. *Transportation Science*, 26(3):161–170, 1992.
- [19] Qing Li, Ming Xu, and Yinhai Wang. Modeling travel time variability with lognormal distribution. *Transportation Research Record*, 2490(1):47–54, 2015.
- [20] Bring Logistics. Customer preferences in last-mile deliveries: Flexible windows and urban density effects. *Industry Report*, 2021. Available via company white papers.

- [21] Mohammadreza Nazari, Afshin Oroojlooy, Lawrence Snyder, and Martin Takáč. Reinforcement learning for solving the vehicle routing problem. In *Proceedings of Advances in Neural Information Processing Systems*, pages 9861–9871, 2018.
- [22] Jorge Oyola, Halvard Arntzen, and David L. Woodruff. The stochastic vehicle routing problem, a literature review, part i: Models. *EURO Journal on Transportation and Logistics*, 7(3):193–221, 2018.
- [23] Jorge Luis Oyola, Halvard Arntzen, and David L Woodruff. The stochastic vehicle routing problem: A literature review, part i: Models. *EURO Journal on Transportation and Logistics*, 7(3):193–221, 2018.
- [24] Nikica Peric, Slaven Begovic, and Vinko Lesic. Adaptive memory procedure for solving real-world vehicle routing problem. *arXiv preprint arXiv:2403.04420*, 2024.
- [25] Nikica Perić, Slaven Begović, and Vinko Lesić. Adaptive memory procedure for solving real-world vehicle routing problem. *arXiv preprint arXiv:2403.04420*, 2024.
- [26] Laurent Perron and Frédéric Didier. Cp-sat.
- [27] David Schrank, Bill Eisele, Tim Lomax, et al. 2021 urban mobility report. *Texas A&M Transportation Institute*, 2021.
- [28] Federal Highway Administration U.S. Department of Transportation. Manual on uniform traffic control devices (mutcd), 2009 edition, 2009. Accident and incident classification and duration guidelines.
- [29] Ron van Duin, Tolga Bektaş, Murat Bektaş, and Tavares Tan. Attended home deliveries: Preferences and behavioral patterns. *Transportation Research Procedia*, 16:30–39, 2016.
- [30] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- [31] Jiongzhi Zheng, Kun He, Jianrong Zhou, Yan Jin, and Chu-Min Li. Reinforced lin-kernighan-helsgaun algorithms for the traveling salesman problems. *Knowledge-Based Systems*, 260:110144, 2023.

## A Open Infrastructure

To ensure reproducibility, extensibility, and accessibility, we release all components of the benchmark openly on GitHub and Hugging Face. This includes the dataset, instance generator, evaluation engine, and baseline implementations. Evaluation instances can be used out of the box, while the modular codebase allows users to integrate new solvers and adapt evaluation scripts.

A public leaderboard on huggingface<sup>2</sup> serves as the central hub for documentation, instance downloads, and leaderboard submissions. Submissions are validated automatically and ranked by total cost, feasibility, and runtime. All data and code are versioned, containerized (Docker-supported), and designed to support future extensions such as new routing scenarios or solver classes.

We welcome community contributions, including new solvers, datasets, and improvements to documentation or evaluation tools. By sharing the infrastructure broadly, we aim to foster collaboration and accelerate progress in realistic stochastic routing research.

### A.1 Reproducibility Requirements

To maintain transparency and enable fair comparison, submissions intended for leaderboard inclusion or academic publication must satisfy several criteria. Solvers must be evaluated on the official benchmark test set, with all hyperparameters, configuration details, and seed values fully documented. Additionally, we encourage open-source releases or detailed methodological descriptions to ensure

---

<sup>2</sup><https://huggingface.co/spaces/ahmedheakl/SVRP-leaderboard>

algorithm reproducibility. Runtime should be measured using the official script or a clearly defined procedure, consistent across all experiments.

These guidelines help uphold reproducibility standards advocated in combinatorial optimization literature [7, 1] and promote meaningful scientific comparisons under controlled, yet realistic, conditions.

## B Baseline Models

**Ant Colony Optimization (ACO).** Routes are constructed by sampling next locations based on pheromone intensity and heuristic proximity. The pheromone matrix is updated as:

$$\tau_{ij} \leftarrow (1 - \rho)\tau_{ij} + \sum_{k=1}^m \Delta\tau_{ij}^{(k)}, \quad \Delta\tau_{ij}^{(k)} = \begin{cases} \frac{Q}{L^{(k)}}, & \text{if } (i, j) \in \text{tour}^{(k)} \\ 0, & \text{otherwise,} \end{cases} \quad (19)$$

where  $\rho = 0.5$ ,  $m = 50$  ants,  $\alpha = 1$ , and  $\beta = 2$ .

**Tabu Search.** Candidate solutions are evaluated using a penalized cost function:

$$f(S) = \text{Cost}(S) + \lambda \cdot \text{Penalty}(S), \quad (20)$$

where  $\lambda$  is adaptively tuned based on violation severity.

**Learning-Based Methods.** The Attention Model is trained to minimize the expected cost:

$$\mathcal{L}(\theta) = \mathbb{E}_{X \sim \mathcal{D}} [\mathbb{E}_{\pi_\theta(a|X)} [L(a|X)]] . \quad (21)$$

POMO uses multiple rollout agents initialized with distinct permutations. Its gradient signal is computed as:

$$\nabla_\theta J(\theta) = \frac{1}{M} \sum_{m=1}^M \sum_t \nabla_\theta \log \pi_\theta(a_t^m | s_t^m) \cdot (R^m - b), \quad (22)$$

where  $M$  is the number of rollouts and  $b$  is a learned baseline for variance reduction.

## C Detailed Solver Performance Breakdowns

Tables 6,7,8,9,10,11 present a comprehensive performance breakdown of various solvers across multiple configurations for Capacitated VRP (CVRP) and Time Window VRP (TWVRP). Each solver, NN+2opt, Tabu Search, ACO, OR-Tools, and RL-based methods (Attention, POMO), is evaluated under different settings including depot configurations (single depot, multi depot, depots equal to cities), problem sizes (ranging from 10 to 1000 customers), and feasibility constraints. Metrics include total cost, CVR (constraint violation rate), feasibility, runtime, and time window violations. Traditional heuristic solvers (NN+2opt, Tabu, ACO) generally yield competitive costs with increasing runtimes as problem size grows, while OR-Tools offers consistent feasibility but with significantly higher runtimes. Reinforcement learning solvers (Attention, POMO) demonstrate exceptionally fast runtimes (in milliseconds), achieving full feasibility across all tested instances, although their cost can vary notably, especially for large-scale problems where some cost inflation is observed (e.g. POMO on 1000-node CVRP). These results highlight trade-offs between solution quality, computational efficiency, and scalability across solver paradigms.

Table 6: NN+2opt - Detailed Performance Breakdown.

Configuration	Size	Cost	CVR	Feas	Runtime	TW Violations
single depot single vehicule sumDemands	10	2290.7	0.0	1.000	0.0	0.00
multi depot	10	2371.8	0.0	1.000	2.0	0.00
single depot single vehicule sumDemands	20	3736.5	0.0	1.000	0.3	0.00
multi depot	20	3662.9	0.0	1.000	3.2	0.00
single depot single vehicule sumDemands	50	4840.4	0.0	1.000	10.5	0.00
multi depot	50	5626.1	0.0	1.000	14.1	0.00
single depot single vehicule sumDemands	100	6841.4	0.0	1.000	31.8	0.00
multi depot	100	7868.2	0.0	1.000	31.3	0.00
single depot single vehicule sumDemands	200	11268.2	0.0	1.000	125.2	0.00
multi depot	200	11479.2	0.0	1.000	135.5	0.00
single depot single vehicule sumDemands	500	16390.0	0.0	1.000	829.5	0.00
multi depot	500	17551.0	0.0	1.000	826.3	0.00
single depot single vehicule sumDemands	1000	25844.3	0.0	1.000	3545.9	0.00
multi depot	1000	25817.4	0.0	1.000	3493.3	0.00
depots equal city	10	4564.6	3.3	0.967	2.0	0.00
single depot	10	4359.0	3.3	0.967	2.0	0.00
depots equal city	20	8192.2	0.0	1.000	5.7	0.00
single depot	20	8347.0	0.0	1.000	5.0	0.00
depots equal city	50	13666.8	0.0	1.000	14.9	0.00
single depot	50	13882.4	0.7	0.993	11.7	0.00
depots equal city	100	38704.2	6.0	0.940	52.2	0.00
single depot	100	30389.4	1.0	0.990	37.8	0.00
depots equal city	200	89937.2	10.1	0.899	145.2	0.00
single depot	200	55400.9	1.0	0.990	167.8	0.00
depots equal city	500	175711.7	7.7	0.923	1318.5	0.00
single depot	500	118279.0	2.2	0.978	929.7	0.00
depots equal city	1000	244956.8	6.1	0.939	3865.2	0.00
single depot	1000	187829.7	2.7	0.973	3911.5	0.00

Table 7: Tabu Search - Detailed Performance Breakdown.

Configuration	Size	Cost	CVR	Feas	Runtime	TW Violations
single depot single vehicle sumDemands	10	2297.2	0.0	1.000	19.4	0.00
multi depot	10	2373.8	0.0	1.000	13.3	0.00
single depot single vehicule sumDemands	20	3776.7	0.0	1.000	47.6	0.00
multi depot	20	3656.4	0.0	1.000	33.8	0.00
single depot single vehicule sumDemands	50	4897.0	0.0	1.000	79.8	0.00
multi depot	50	5749.3	0.0	1.000	102.9	0.00
single depot single vehicule sumDemands	100	6981.9	0.0	1.000	170.0	0.00
multi depot	100	8058.6	0.0	1.000	169.2	0.00
single depot single vehicule sumDemands	200	11417.8	0.0	1.000	373.9	0.00
multi depot	200	11602.8	0.0	1.000	314.2	0.00
single depot single vehicule sumDemands	500	16554.8	0.0	1.000	1270.4	0.00
multi depot	500	17676.2	0.0	1.000	1445.1	0.00
single depot single vehicule sumDemands	1000	25995.4	0.0	1.000	4647.9	0.00
multi depot	1000	25879.7	0.0	1.000	4544.5	0.00
depots equal city	10	3966.1	3.3	0.667	185.9	0.00
single depot	10	4067.6	3.3	0.667	193.6	0.00
depots equal city	20	8156.1	0.0	1.000	479.8	0.00
single depot	20	7661.3	0.0	1.000	489.9	0.00
depots equal city	50	13918.7	0.0	1.000	719.3	0.00
single depot	50	14269.3	0.7	0.667	654.4	0.00
depots equal city	100	39031.2	6.0	0.000	2013.6	0.00
single depot	100	30820.4	1.0	0.333	1998.3	0.00
depots equal city	200	90028.5	10.1	0.000	2662.6	0.00
single depot	200	55596.2	1.0	0.000	3014.1	0.00
depots equal city	500	176001.3	8.1	0.000	13851.1	0.00
single depot	500	118726.0	2.2	0.000	11822.7	0.00
depots equal city	1000	244953.3	6.2	0.000	50402.1	0.00
single depot	1000	187945.6	2.7	0.000	42673.2	0.00

Table 8: ACO - Detailed Performance Breakdown.

Configuration	Size	Cost	CVR	Feas	Runtime	TW Violations
single depot single vehicule sumDemands	10	2183.6	0.0	1.000	14.3	0.00
multi depot	10	2325.4	0.0	1.000	11.9	0.00
single depot single vehicule sumDemands	20	3725.9	0.0	1.000	34.6	0.00
multi depot	20	3644.2	0.0	1.000	31.4	0.00
single depot single vehicule sumDemands	50	4840.5	0.0	1.000	165.2	0.00
multi depot	50	5626.2	0.0	1.000	179.5	0.00
single depot single vehicule sumDemands	100	6840.4	0.0	1.000	698.1	0.00
multi depot	100	7868.4	0.0	1.000	678.2	0.00
single depot single vehicule sumDemands	200	11264.3	0.0	1.000	2295.7	0.00
multi depot	200	11473.0	0.0	1.000	2380.3	0.00
single depot single vehicule sumDemands	500	16389.2	0.0	1.000	15573.5	0.00
multi depot	500	17551.6	0.0	1.000	16468.6	0.00
single depot single vehicule sumDemands	1000	25840.7	0.0	1.000	58364.4	0.00
multi depot	1000	25815.8	0.0	1.000	59341.2	0.00
depots equal city	10	3931.6	3.3	0.667	9.4	0.00
single depot	10	3819.2	3.3	0.667	9.6	0.00
depots equal city	20	7714.2	0.0	1.000	34.2	0.00
single depot	20	7749.4	0.0	1.000	34.1	0.00
depots equal city	50	13535.4	0.0	1.000	166.9	0.00
single depot	50	13872.4	0.7	0.667	143.6	0.00
depots equal city	100	37800.2	6.0	0.000	629.4	0.00
single depot	100	30389.5	1.0	0.333	679.0	0.00
depots equal city	200	89937.2	10.1	0.000	2556.8	0.00
single depot	200	55401.8	1.0	0.000	2327.0	0.00
depots equal city	500	175711.1	7.7	0.000	15299.3	0.00
single depot	500	118280.2	2.2	0.000	14781.5	0.00
depots equal city	1000	244999.0	6.1	0.000	70932.6	0.00
single depot	1000	187332.2	2.8	0.000	54846.8	0.00

Table 9: OR-Tools - Detailed Performance Breakdown.

Configuration	Size	Cost	CVR	Feas	Runtime	TW Violations
single depot single vehicule sumDemands	10	2049.2	0.0	1.000	1037.9	0.00
multi depot	10	2167.6	0.0	1.000	1003.3	0.00
single depot single vehicule sumDemands	20	3238.9	0.0	1.000	999.5	0.00
multi depot	20	3142.2	0.0	1.000	1002.6	0.00
single depot single vehicule sumDemands	50	3773.4	0.0	1.000	1015.9	0.00
multi depot	50	4714.2	0.0	1.000	1015.9	0.00
single depot single vehicule sumDemands	100	6283.5	0.0	1.000	1046.5	0.00
multi depot	100	6250.4	0.0	1.000	1048.8	0.00
single depot single vehicule sumDemands	200	9198.8	0.0	1.000	1174.7	0.00
multi depot	200	8956.2	0.0	1.000	1185.4	0.00
single depot single vehicule sumDemands	500	15677.5	0.0	1.000	2129.5	0.00
multi depot	500	15883.2	0.0	1.000	2085.2	0.00
single depot single vehicule sumDemands	1000	25844.3	0.0	1.000	8412.4	0.00
multi depot	1000	25816.3	0.0	1.000	9434.5	0.00
depots equal city	10	4564.7	3.3	0.967	12.6	0.00
single depot	10	4359.0	3.3	0.967	3.6	0.00
depots equal city	20	8192.3	0.0	1.000	8.2	0.00
single depot	20	8346.9	0.0	1.000	7.2	0.00
depots equal city	50	13666.7	0.0	1.000	30.3	0.00
single depot	50	13882.3	0.7	0.993	27.7	0.00
depots equal city	100	38704.1	6.0	0.940	108.6	0.00
single depot	100	30389.3	1.0	0.990	87.8	0.00
depots equal city	200	89937.5	10.1	0.899	345.3	0.00
single depot	200	55401.8	1.0	0.990	329.8	0.00
depots equal city	500	175711.4	7.7	0.923	2010.0	0.00
single depot	500	118279.4	2.2	0.978	2020.5	0.00
depots equal city	1000	244998.0	6.1	0.939	8273.1	0.00
single depot	1000	187830.1	2.7	0.973	8464.4	0.00



Table 10: RL Algorithms – Detailed Performance on CVRP (runtimes in ms).

Solver	Configuration	Size	Cost	CVR	Feas	Runtime (ms)	TW Violations
Attention	single depot single vehicule sumDemands	10	2364.12	0.00	1.000	0.365	0.00
POMO	single depot single vehicule sumDemands	10	2312.68	0.00	1.000	0.282	0.00
Attention	single depot single vehicule sumDemands	20	3222.68	0.00	1.000	0.269	0.00
POMO	single depot single vehicule sumDemands	20	3341.56	0.00	1.000	0.279	0.00
Attention	single depot single vehicule sumDemands	50	5803.63	0.00	1.000	0.304	0.00
POMO	single depot single vehicule sumDemands	50	5920.19	0.00	1.000	0.287	0.00
Attention	single depot single vehicule sumDemands	100	8553.26	0.00	1.000	0.319	0.00
POMO	single depot single vehicule sumDemands	100	16983.50	0.00	1.000	0.319	0.00
Attention	single depot single vehicule sumDemands	200	13228.84	0.00	1.000	0.353	0.00
POMO	single depot single vehicule sumDemands	200	12726.96	0.00	1.000	0.360	0.00
Attention	single depot single vehicule sumDemands	500	22496.94	0.00	1.000	0.463	0.00
POMO	single depot single vehicule sumDemands	500	88789.44	0.00	1.000	0.506	0.00
Attention	single depot single vehicule sumDemands	1000	37430.47	0.00	1.000	0.649	0.00
POMO	single depot single vehicule sumDemands	1000	184656.10	0.00	1.000	0.689	0.00

Table 11: RL Algorithms – Detailed Performance on TWVRP (runtimes in ms).

Solver	Configuration	Size	Cost	CVR	Feas	Runtime (ms)	TW Violations
Attention	single depot	10	3 940.38	0.00	1.000	0.916	0.00
POMO	single depot	10	3 854.6	0.00	1.000	0.707	0.00
Attention	single depot	20	6 504.73	0.00	1.000	1.780	0.00
POMO	single depot	20	6 744.7	0.00	1.000	1.841	0.00
Attention	single depot	50	29 132.94	0.00	1.000	0.731	0.00
POMO	single depot	50	29 718.0	0.00	1.000	0.689	0.00
Attention	single depot	100	57 778.84	0.00	1.000	0.864	0.00
POMO	single depot	100	114 726.7	0.00	1.000	0.864	0.00
Attention	single depot	200	113 742.27	0.00	1.000	0.868	0.00
POMO	single depot	200	109 427.1	0.00	1.000	0.886	0.00
Attention	single depot	500	271 201.60	0.00	1.000	1.412	0.00
POMO	single depot	500	438 502.6	0.00	1.000	1.412	0.00
Attention	single depot	1000	531 470.88	0.00	1.000	1.638	0.00
POMO	single depot	1000	611 307.8	0.00	1.000	1.672	0.00

## C.1 Qualitative Results

As shown in figures 5, 6, 7, 8, 9, 10, 11, 12, 13, and 14, we qualitatively observe that for CVRP instances with a small number of customers, both Attention and POMO models, as well as classical methods (ACO, NN2OPT, and OR-Tools), generate highly structured and near-optimal routes. As the number of customers increases, route complexity grows, making it harder for models to preserve efficiency and structure. For TWVRP, the models' priority shifts toward satisfying delivery time windows, often at the expense of distance optimization. This results in routes that appear less spatially coherent but better aligned with temporal constraints.

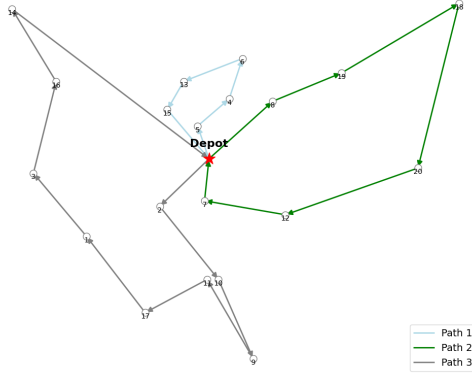


Figure 5: CVRP 20 customers – Attention Model

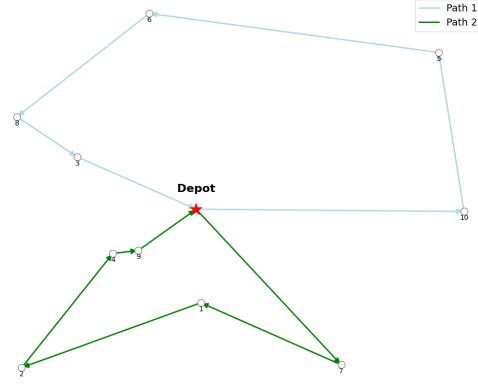


Figure 6: CVRP 10 customers – POMO

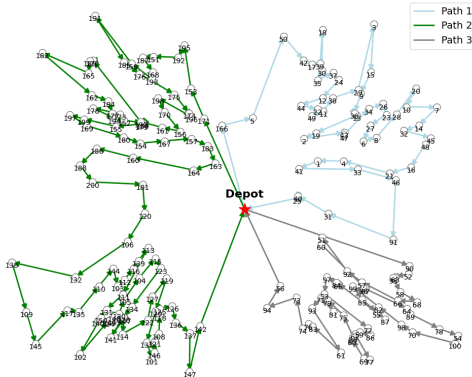


Figure 7: CVRP 200 customers – Attention Model

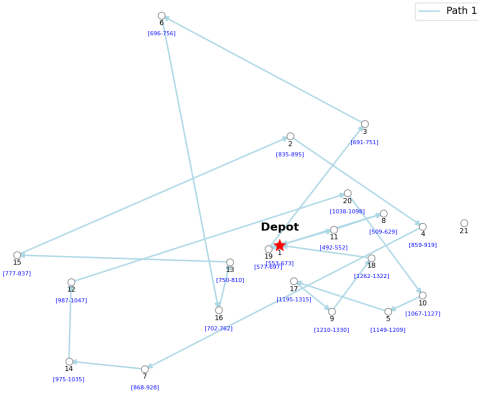


Figure 8: TWVRP 20 customers – Attention Model

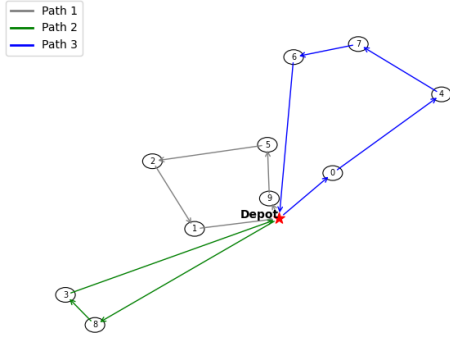


Figure 9: CVRP 10 customers – ACO

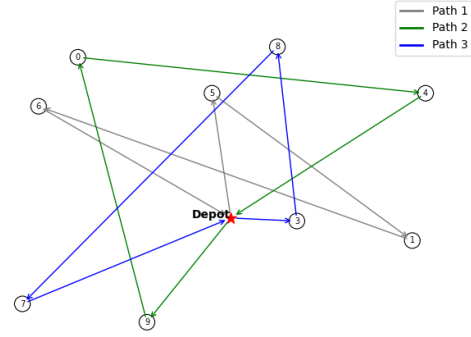


Figure 10: TWVRP 10 customers – ACO

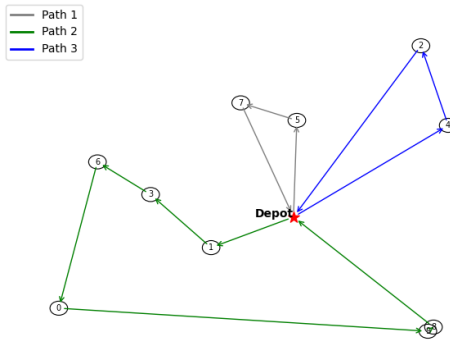


Figure 11: CVRP 10 customers – NN2OPT

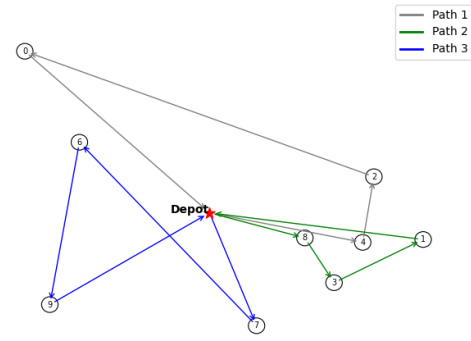


Figure 12: TWVRP 10 customers – NN2OPT

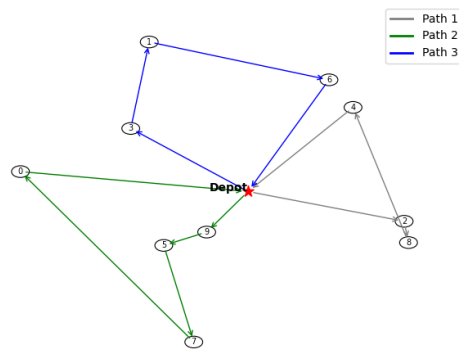


Figure 13: CVRP 10 customers – OR-Tools

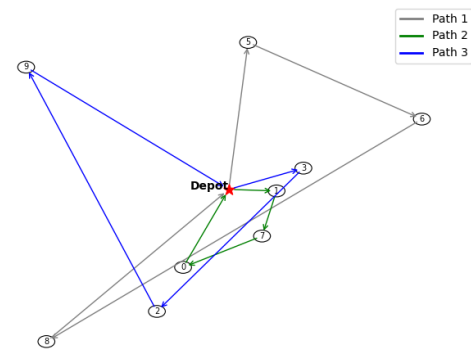


Figure 14: TWVRP 10 customers – OR-Tools

## D Reinforcement Learning

### D.1 Problem Formulation

We model both the Capacitated Vehicle Routing Problem (CVRP) and Vehicle Routing Problem with Time Windows (VRPTW) as a Markov Decision Process (MDP)  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$ , where each state  $s_t \in \mathcal{S}$  encodes the vehicle’s current position, remaining capacity, visited set (and only for VRPTW the current time and per-customer time windows  $[e_i, \ell_i]$ ). Actions  $a_t \in \mathcal{A}(s_t)$  select the next customer, and transitions  $P(s_{t+1} | s_t, a_t)$  deterministically update the tour while, in VRPTW, adding stochastic delays.

The reward is  $r(s_t, a_t) = -d_{i,j} - \tau [t_{\text{arrive}} > \ell_i]$  when visiting customer  $j$ , with  $d_{i,j}$  the Euclidean distance and  $\tau$  a large penalty for time-window violations, and zero upon return to the depot. We follow a constructive, autoregressive decoding: at each step we append one customer until all are visited.

### D.2 Policy

We adopt the encoder–decoder with multi-head attention of Kool [16]. Given embedded node features  $\mathbf{x}_i \in \mathbb{R}^d$ , each of the  $L$  encoder layers applies multi-head self-attention. At step  $t$ , with context embedding  $\mathbf{h}_t$ , we score each remaining node  $j$  by  $u_{t,j} = \mathbf{v}^\top \tanh(W_1 \mathbf{h}_t + W_2 \mathbf{x}_j)$  and define  $\pi_\theta(a_t = j | s_t) = \exp(u_{t,j}) / \sum_{k \notin \mathcal{V}_t} \exp(u_{t,k})$ .

We optimize the policy by maximizing the expected return  $J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[R(\tau)]$  using two constructive, autoregressive policy-gradient methods. A constructive policy builds a complete solution by sequentially selecting one customer at a time until the tour is finished, while an autoregressive policy conditions each action on the history of previous choices, enabling the network to capture dependencies across steps.

We first apply REINFORCE [30], which updates parameters via  $\nabla_\theta J(\theta) = \mathbb{E}[\sum_t \nabla_\theta \log \pi_\theta(a_t | s_t) (R(\tau) - b(s_t))]$ , where  $b(s_t)$  is a rollout baseline obtained by greedy decoding; then POMO [17] samples  $K$  different start nodes per instance, computes returns  $R_k$  and a shared baseline  $\bar{R} = \frac{1}{K} \sum_k R_k$ , and applies  $\nabla_\theta J(\theta) = \frac{1}{K} \sum_{k=1}^K \nabla_\theta \log \pi_\theta(\tau_k) (R_k - \bar{R})$ . REINFORCE offers simplicity and unbiased gradients, while POMO’s shared baseline exploits VRP permutation symmetry for variance reduction; together they provide a strong comparison between a classical Monte Carlo approach and a state-of-the-art, variance-reduced VRP-specific algorithm.

### D.3 Training Details

All models were implemented in the RL4CO framework and trained end-to-end with Adam at a learning rate of  $10^{-4}$ . For CVRP with REINFORCE we used a batch size of 512 and generated 100 000 synthetic instances on the fly; for VRPTW with POMO we used batch size 64 and 1 000 000 instances. Validation employed greedy decoding under nominal travel-time conditions. VRPTW environments included log-normal delays calibrated to traffic data, Gaussian time-of-day kernels, and Poisson accident events, with infeasible actions heavily penalized to enforce time windows.

### D.4 Evaluation on SVRPBench

After training, we converted each of the 500+ SVRPBench instances into the RL4CO environment format and ran the trained policies in greedy mode, selecting at each step  $a_t = \arg \max_j \pi_\theta(a_t = j | s_t)$ . To assess robustness, we then simulated each resulting tour under multiple sampled delay realizations and reported average tour length and feasibility rates. Despite domain shift, attention-based RL policies maintained high feasibility and near-optimal costs across all problem sizes.