# Guidelines for the Computational Testing of Machine Learning approaches to Vehicle Routing Problems

Luca Accorsi[1], Andrea Lodi[2, 3], and Daniele Vigo[1, 4]

[1]*Department of Electrical, Electronic and Information Engineering "G. Marconi", University of Bologna, Italy*
[2]*Canada Excellence Research Chair in Data Science for Decision Making, École Polytechnique de Montréal, Canada*
[3]*Mila, Quebec Artificial Intelligence Institute, Canada*
[4]*CIRI ICT, University of Bologna, Italy*

`{luca.accorsi4,daniele.vigo}@unibo.it`
`andrea.lodi@polymtl.ca`

**Abstract**

Despite the extensive research efforts and the remarkable results obtained on Vehicle Routing Problems (VRP) by using algorithms proposed by the Machine Learning community that are partially or entirely based on data-driven analysis, most of these approaches are still seldom employed by the Operations Research (OR) community. Among the possible causes, we believe, the different approach to the computational evaluation of the proposed methods may play a major role. With the current work, we want to highlight a number of challenges (and possible ways to handle them) arising during the computational studies of heuristic approaches to VRPs that, if appropriately addressed, may produce a computational study having the characteristics of those presented in OR papers, thus hopefully promoting the collaboration between the two communities.

## 1 Introduction

Recently, several attempts to use machine learning (ML) to deal with combinatorial optimization problems (COPs) have been proposed. Indeed, ML promises to support the design of algorithms as well as the resolution process, by decreasing the need for hand-crafted and specialized solution approaches, which are notably known to require high expertise and a huge time investment for being developed. The incredible advances in deep learning (DP, Goodfellow, Bengio, and Courville (2016)) coupled with increasingly powerful hardware, have led to very promising results for a variety of COPs (Bengio, Lodi, and Prouvost (2020) and Mazyavkina et al. (2020)).

In this work, we focus on COPs arising in the area of vehicle routing. Vehicle routing problems (VRPs, Toth and Vigo (2014)) are, in fact, increasingly receiving attention in the ML community both because of their relevance in real-world applications and the computational challenges they pose (Vesselinova et al. (2020)). Indeed, VRPs have been the test-bed for novel neural network architectures such as pointer networks (Vinyals, Fortunato, and Jaitly (2015)), graph embedding networks (H. Dai, B. Dai, and Song (2016) and Khalil et al. (2017)) and attention-based models (Kool, Hoof, and Welling (2019) and Deudon et al. (2018)) achieving stunning results on a variety of problems. Currently, most of the proposed approaches aim at learning constructive heuristics, which sequentially extend a partial solution, possibly employing additional procedures such as sampling and beam search (see e.g., Bello et al. (2017) and Hottung, Kwon, and Tierney (2021)). Few others, such as Wu et al. (2020) and Chen and Tian (2019), instead,

focus on learning improvement heuristics to guide the exploration of the search space and iteratively refine an existing solution.

Despite of the relevance of the problems, the increasing attention they are receiving in the ML community and the remarkable results achieved so far, these techniques have not yet become widespread in the Operations Research (OR) community. This may be because of their novelty and the nontrivial effort required to adopt them by the OR community that typically has a different background. But also, at the same time, the computational testing provided by some of the proposed ML methods may not be very convincing with respect to the standard practices in the OR and VRP community. The current work focuses on this second aspect by highlighting important points to consider during the computational testing and providing guidelines and methodologies that we believe are relevant from an OR perspective. The remainder of the paper is organized as follows. Section 2 surveys crucial aspects faced during a computational testing such as the selection of proper benchmark instances and baseline algorithms, as well as techniques to properly compare different solution approaches. Section 3 provides concrete examples and pointers to representative computational studies found in recent papers. Finally, concluding remarks are given in Section 4.

## 2 Guidelines for a Computational Testing

In this section, we highlight important points an OR practitioner would consider crucial when examining the computational results section of a novel approach.

### 2.1 Benchmark Instances and Problem Definition

The training phase of a ML-based approach typically requires a large number of VRP instances sharing the same characteristics. A common way to generate these instances is by randomly defining their characteristics, sampled from arbitrarily defined distributions. In the following, we argue that adopting a common problem description as well as a common set of benchmark instances is extremely important to correctly assess the potentialities of a novel approach.

*Problem description and objectives.* The term VRP identifies a class of problems containing an enormous number of different variants. Among the most studied ones, we have the Capacitated Vehicle Routing Problem (CVRP), the Vehicle Routing Problem with Time Windows (VRPTW), and the Vehicle Routing Problem with Pickup and Deliveries. In the following, if not specified differently, we limit our treatment to the CVRP which is often taken as a reference problem in ML-based approaches and, despite its simple definition, is still a challenging problem for both traditional and innovative approaches. We refer to Chapter 1 of Toth and Vigo (2014) for a thorough overview of VRPs. Specific VRPs have widely accepted formulations and precise problem definitions. For example, modern CVRP instances do not fix a-priori the number of routes a solution should have or heuristic approaches to the VRPTW typically consider a hierarchical objective: first minimize the number of vehicles and then the routing cost. It is thus clear that the specific definition of the problem has a crucial impact on the obtained results.

*Test instances representativeness.* Despite VRPs being $\mathcal{NP}$-hard, the actual challenge posed by a specific instance is highly dependant on several factors such as customers and depot location, the vehicle capacity and the customers demand distribution. The VRP community has thus identified, for each specific problem in the VRP class, a set of benchmark instances that is currently considered to be relevant for testing modern approaches. Thus, in addition to possible instances defined by the ML community, we highlight the importance of using, whenever possible, instances derived from these widely recognized, used and studied datasets.

As an example, in Uchoa et al. (2017), the authors introduced the $\mathbb{X}$ instances for the CVRP, thoroughly describing their generation process. This process can then be emulated to generate (possibly smaller-sized) more representative CVRP instances, as it was done in Kool, Hoof, Gromicho, et al. (2021) and in Hottung, Kwon, and Tierney (2021) (appendix). In Wu et al. (2020), the authors directly used a subset of $\mathbb{X}$ instances along with distributions more commonly used in other ML works.

As an additional example, consider the VRPTW for which the so-called Solomon instances and their extension proposed by Solomon (1987) and Gehring and Homberger (1999), respectively, are the current benchmark. Also for them, Solomon (1987) describes the procedure used to define the time window constraints that can thus be considered when defining new instances.

*Repositories of instances.* Together with papers introducing or using certain sets of instances (whose authors could be contacted to retrieve), we mention two among the most popular repositories of VRP instances. Namely, VRP-REP collects instances for more than 50 different VRP variants, best known solution values, and references to papers obtaining these results. In addition, CVRPLIB contains instances and up-to-date best-known solutions of CVRP instances. Finally, small instances commonly used in the past or in exact methods can be found in TSPLIB and OR-LIBRARY.

## 2.2 Baseline Algorithms

The selection of a proper baseline algorithm is extremely important, failing in this task would hinder the objective evaluation of the potential of a novel approach. Indeed, since results (computing time, solution quality, and possibly more sophisticated measures as detailed in Section 2.3) are crucial to compare different solution approaches over a common set of problems and instances, a wrong baseline may distort their interpretation, undermining the whole validation process. Despite the purpose of ML-based approaches not being that of outperforming highly specialized solvers, but rather that of proposing versatile tools not requiring high-levels of manual engineering, the comparison should still occur against the best performing algorithms to better comprehend the tradeoff between data-driven and ad-hoc algorithms.

*Include the best available algorithms.* Along with simple baselines and competing ML-based methods, we argue that one should consider the inclusion of the best available algorithms proposed by the OR community for each specific VRP.

The selection of baseline algorithms is often guided by the availability of free-to-use or open-source reliable software packages. Fortunately, more and more researchers are publishing the source code of heuristic as well as exact state-of-the-art VRP solvers that can be freely used for research activities.

**Heuristic solvers**   The widely used Google OR-Tools (Perron and Furnon (2019)) is erroneously considered by most ML papers to be among the best open-source VRP solvers (see, e.g., Nazari et al. (2018)) while achieving on the CVRP far-from-optimal results on the $\mathbb{X}$ instances (see Vidal (2020)). Much better open-source solvers for the CVRP are fortunately available. Along with LKH-3 (Helsgaun (2017)), which is already widely used by the ML community for solving VRPs, we mention HGS-CVRP (Vidal (2020)) and FILO (Accorsi and Vigo (2020)) as highly effective and efficient open-source heuristic solvers for the CVRP that, on the widely studied $\mathbb{X}$ instances of the CVRP (Uchoa et al. (2017)), produce superior results compared with LKH-3 (see Cavaliere, Bendotti, and Fischetti (2020)). Finally, we mention SISR (Christiaens and Vanden Berghe (2020)), that, despite not being already available in terms of source code, is conceptually simple and easy to implement, yet providing state-of-the-art results on a great variety of VRPs.

**Exact solvers** Several papers solve small instances to optimality by using general-purpose optimization solvers such as Gurobi or CPLEX. Despite the noble attempt, trying to directly solve simple compact VRP formulations by using a generic branch-and-cut approach would soon turn out to be an extremely challenging task. Indeed, several papers report the ability of solving only very small instances (e.g., CVRPs with about 20 customers). Instead, VRPSolver (Pessoa et al. (2020)), a freely available (for academic purposes) exact solver specialized for routing problems, should be considered for serious and reliable testing of VRPs. Indeed, VRPSolver combines a branch-cut-and-price algorithm with other sophisticated techniques specifically designed for VRPs such as route enumeration, state space relaxation, and others being able to consistently solve CVRP instances with up to 200 customers (a size already out of reach for most ML-based approaches proposed so far).

## 2.3 Algorithms Comparison

Comparing algorithms having a completely different nature is an extremely challenging task. On the one hand, traditional solution approaches proposed by the OR community are designed to handle set of instances having a broad range of different characteristics. Moreover, these approaches are typically tuned to achieve, with a single set of parameter values, results that are on average good over all tested instances. On the other hand, ML-based approaches generally need to treat every instance distribution separately, requiring a specific tuning and thus additional training (possibly taking up to several weeks of computing time). In addition, traditional OR algorithms are (almost) always executed on CPUs using a single thread, while ML-based approaches naturally benefit from running on massively parallel hardware architectures such as GPUs. Finally, the programming language may also play a role. In fact, traditional OR algorithms are usually implemented in highly efficient languages such as C++, whereas ML-based approaches typically use Python that mixes slow interpreted code with efficient native libraries.

*Facilitate comparisons (i.e., run all solvers by using a common configuration).* In production, algorithms should obviously make fully use of the best existing available technologies. In fact, even traditional algorithms may contain (portions of) embarrassingly parallel code that would thus benefit from being run on multiple threads. As an example, a common dynamic programming procedure employed in VRP exact solvers would greatly benefit from a GPU implementation compared to a traditional sequential version (Boschetti, Maniezzo, and Strappaveccia (2017) report speedup of up to 40 times). Another well-known example is the parallel implementation of Branch & Bound algorithms (see, e.g., Crainic, Le Cun, and Roucairol (2006)). Despite the clear time-savings of a parallel implementation, experimental evaluation asks to reduce at a minimum the different factors that would render comparisons among approaches unnecessarily more challenging. It is thus commonly accepted to consider single-threaded algorithms run on standard CPU architectures. A very interesting information, that would promote a direct comparison of newly proposed algorithms with existing ones, consists in including also the computing time taken by the algorithm when run with the above defined settings. If running all the experiments both on GPU and on CPU with a single thread would be computationally prohibitive, we argue that one should consider adding a measure of the speedup associated with the model inference when run on a GPU rather than on a CPU together with the total algorithmic time (in which everything except the model inference is run on a CPU with a single thread), and the fraction of time spent doing inference on GPU. This way, the total running time of the algorithm on a CPU with a single thread could be easily estimated. Another possibility sees the inclusion of a rough measure of the number of CPU cores needed to match the GPU capacity as it was done in Hottung, Kwon, and Tierney (2021).

*Convert computing time to a common scale.* When using results published on other papers, because the source code may not be available or running again the experiments may be too time consuming, a com-

mon practice consists in roughly scaling the computing time to a common base by using appropriate factors. A practice increasingly adopted in the VRP community consists in using the single-thread rating defined by PassMark. So that, given a base processor, say an Intel Xeon CPU E3-1245 v5 having a single-thread rating of 2277, and a target processor, say an Intel Core2 Duo T5500 having a single-thread rating of 594, the computing time of an algorithm run on the target processor is reduced of a factor $\approx 3.83$. The comparison is still rough, being the CPU just one among the several components affecting the overall performance. However, this is one of the possible approaches generally accepted as sound by the VRP community. Note that, the above considerations assume all algorithms have been run on an unloaded system and on a CPU by using a single thread.

*On the comparison with exact solvers.* Exact solvers are often used as baseline algorithms when approaching small-to-medium size instances. When reporting the computing time spent by an exact solver, especially if its results are compared with a heuristic algorithm, it should be considered that the former may find very good quality solutions early during the run and then spend the majority of the time to prove optimality. Thus, when reporting results obtained by an exact solver, an additional column could be added showing the computing time in which the last solution (i.e., the optimal one, if the solver is not prematurely stopped) was found. Despite being true that the solver does not have a termination criterion to know whether the solution at hand is optimal or not, being compared with heuristics, this approach would still provide a feeling on the convergence speed of the solver. Another possibility, which however requires a much more detailed collection and processing of data, consists in realizing a convergence profile chart (see later). We finally mention the primal integral introduced by Achterberg, Berthold, and Hendel (2012), as a measure able to take into account the overall solution process in terms of convergence towards the optimal (or best known) solution over the entire solving time.

*Consider the statistical relevance of the results.* In the past, it was common practice to use just average percentage errors with respect to the best known solution to compare the performance of different algorithms. However, more recently, the inclusion of simple statistical tests to objectively assess the differences among algorithms is increasingly becoming popular in the VRP community. A common practice consists in using a one-tailed Wilcoxon signed-rank test (see Wilcoxon (1945)) possibly coupled with correction methods when multiple comparisons involving the same data are performed (e.g., the Bonferroni correction, see Dunn (1961)). These tests are used to determine whether two sets of paired observations, for which no assumption can be done on their distribution, are statistically different, and thus, whether two algorithms are considered to provide equivalent results. A common tool for performing these statistical tests is the R language (R Core Team), which allows to execute them with just a few lines of code. An example of a test application can be found in Section 3.3.

*Always compare averages.* When experimenting with algorithms containing randomized components, we argue that one should make comparison by considering the average results obtained over a reasonable number of runs. Typically, the used number of runs are 10 or 50, that, despite not being statistically significant, allow to qualitatively identify whether an algorithm provides stable or highly variable results.

*Include charts.* Charts allow to visually compare several algorithms at once, thus providing a fast and effective way to view the results, that, when coupled with tables and statistical tests, give a satisfying and thorough picture of the computational studies for a set of instances. Among the most used charts, we mention

- the performance chart, relating the average normalized computing time with the average solution quality (e.g., the percentage gap with respect to a reference value, which is typically that of the optimal or best known solution) of each algorithm in the comparison, see Figure 1 (left);
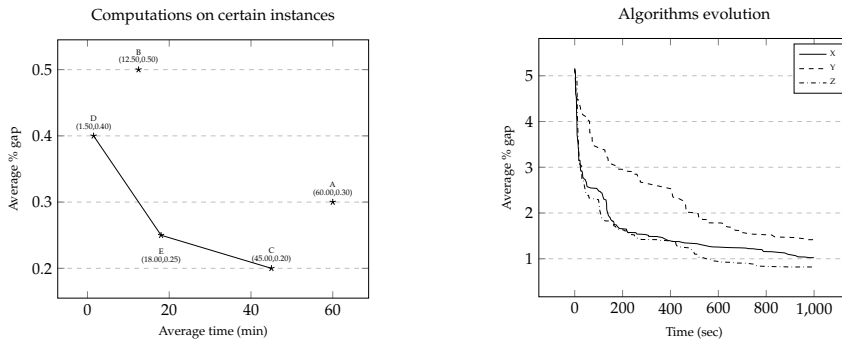
5

Figure 1: On the left side, the performance chart showing for algorithms A, B, C, D, and E the average normalized time *t* and solution quality obtained in a number *n* of runs with different seeds. Algorithms C, D, and E dominate A and B. On the right side, the search trajectory defined by the average gap found at a given time instant for algorithms X, Y and Z.

- the convergence profile chart, showing for each time instant the average solution quality for the compared algorithms, see Figure 1 (right).

The performance chart clearly identifies Pareto optimal algorithms and dominated ones, whereas the convergence profile chart shows their converge speed when executed for a specific period of time. Convergence profile charts can be used for improvement heuristics but also for simpler constructive heuristics provided that they include additional iterative procedures used to further improve the solution quality (e.g., the active search of Bello et al. (2017)).

*On the choice of the programming language.* (Baseline) Algorithms should be implemented efficiently. This may include using appropriate data structures as well as programming languages producing directly executable machine code. An efficient implementation of a competing algorithm should not be considered negatively.

*Consider the analysis of the various algorithm components.* Several papers include additional analysis on the components of the proposed algorithm. This may include the behavior of the algorithm when some parameters are changed as well as the contribution of individual components to the overall final results (e.g., the average improvement of different local search operators analyzed throughout the algorithm execution). This kind of analysis is both considerably appreciated in the VRP community and extremely important to grasp insights on the overall contribution and usefulness of the different components of an algorithm.

## 3   Examples of a Computational Studies

To make the above suggestions more concrete, in this section, we report extracts of and pointers to representative computational studies found in recent papers on the CVRP. We will consider two scenarios: in the first one, we assume to have all source codes available, whereas in the second one, we assume at least one of the competing algorithm source code is not available to the tester.

### 3.1   Scenario 1: all algorithms source codes are available

This is the simplest and most fortunate setting. Indeed, given enough time and processing power, all algorithms can be run on exactly the same platform and for the same amount of time. In this context the

| Instance | BKS | HGS-CVRP | | SISR | | OR-Tools | |
|---|---|---|---|---|---|---|---|
| | | Avg | Gap | Avg | Gap | Avg | Gap |
| X-n101-k25 | 27591 | 27591.0 | 0.00 | 27593.3 | 0.01 | 27977.2 | 1.40 |
| X-n106-k14 | 26362 | 26381.4 | 0.07 | 26380.9 | 0.07 | 26757.5 | 1.50 |
| X-n110-k13 | 14971 | 14971.9 | 0.00 | 14972.1 | 0.01 | 15099.8 | 0.86 |
| ... | | | | | | | |
| X-n979-k58 | 118987 | 119247.5 | 0.22 | 119108.2 | 0.10 | 123885.2 | 4.12 |
| X-n1001-k43 | 72359 | 72748 | 0.54 | 72533.1 | 0.24 | 78084.7 | 7.91 |
| Mean | | | 0.11 | | 0.19 | | 4.01 |

Table 1: Minimal table showing the solution quality obtained by algorithms run for the same amount of time on a common platform.

convergence profile chart perfectly shows the evolution of each algorithm, making evident its speed to reach certain quality levels.

Table 1 reports a rearranged extract of Tables 1-3 proposed in Vidal (2020) comparing three of the above mentioned algorithms (more specifically HGS-CVRP, SISR, and OR-Tools) run on a common platform for the same computing time and over the same $\mathbb{X}$ instances of the CVRP. In particular, the table shows for each competing algorithm the average solution value (Avg) and the associated average gap (Gap) computed considering a certain number of runs of the algorithm when it includes randomized components. The gap is computed with respect to a reference best known solution value (BKS) as Gap = $100 \cdot (\text{Avg} - \text{BKS})/\text{BKS}$. Additional useful columns could include the worst and best gap obtained, so as to better examine the variability of the quality for the algorithm on the dataset. Finally, we refer to Figures 3 and 4 of Vidal (2020) for examples of convergence profile charts for the above-mentioned algorithms.

## 3.2 Scenario 2: at least one algorithm source code is not available

This scenario, which is still unfortunately frequent in the VRP community, makes the comparison of results much more challenging. The common case sees only the presence of tables showing the performance of a proposed algorithm over a set of benchmark instances when run on a certain platform. First of all, an assumption is required when reviewing these tables reporting the computational results. In particular, we have to assume that the results published in the paper, which are inevitably obtained with a specific tuning of parameters, are the values on which the authors desire to compare with other approaches. We shall note that these parameters do include the termination criterion, which was then selected as a design choice, and considered to provide competitive results (otherwise a different criterion would have been selected). Since data on the convergence of the algorithm are typically not available, the comparison can then be performed over the published results by first normalizing the computing time in the best possible way and then analyzing the bi-objective perspective provided by the performance chart shown in 1 (left) showing non dominated algorithms for a fixed configuration of their parameters.

## 3.3 Statistical validation of the results

In both the above scenarios (and provided that the average solution quality obtained by an algorithm is available for each instance in the dataset), the (final) solution quality could be assessed with simple statistical tests. This procedure is useful especially when the proposed methods do not clearly dominate the others, for example because they obtain similar average percentage gaps on the considered benchmark instances.

As an example, in the following we compare HGS-CVRP, SISR, and OR-Tools on the $\mathbb{X}$ instances. Data have been taken from Tables 1 and 3 of Vidal (2020) and is summarized in the boxplots of Figure 2.

Similarly to what was done in Christiaens and Vanden Berghe (2020), we can assess the results obtained by our proposed algorithm, say HGS-CVRP, against the remaining ones, by conducting a one-tailed Wilcoxon signed-rank test in which we consider a null hypothesis $H_0$

$$H_0 : \textsc{AvgSolCost}(\text{HGS-CVRP}) = \textsc{AvgSolCost}(\text{X}),$$

and an alternative hypothesis $H_1$

$$H_1 : \textsc{AvgSolCost}(\text{HGS-CVRP}) > \textsc{AvgSolCost}(\text{X}),$$

where $X$ can be SISR and OR-Tools. A hypothesis is rejected when its $p$-value is lower than a significance level $\alpha$. In particular, we have that

- failing to reject $H_0$ means that the average results of the two methods are not statistically different;

- whereas, when $H_0$ is not rejected, the average results are statistically different and the alternative hypothesis $H_1$ can be tested to find whether they are greater than those of a competing method. Rejecting $H_1$ thus implies that HGS-CVRP performs better than the competing method.

Moreover, as mentioned in Section 2.3, when performing multiple comparisons involving the same data, the probability of erroneously rejecting a null hypothesis increases. To control these errors, the significance level $\alpha$ is adjusted to lower values. Bonferroni correction (Dunn (1961)) is a simple method that can be used for this purpose. In particular, given $n$ comparisons, the significance level is set to $\alpha/n$.

In the above comparison we tested a total number of $n = 2$ hypothesis corresponding to the partitioning of instances (1, all $\mathbb{X}$ instances together) and to the two hypothesis (2, $H_0$ and $H_1$). Assuming an initial significance level $\alpha_0 = 0.025$, the adjusted value becomes $\alpha = 0.025/2 = 0.0125$.

We can compute the $p$-values, which are shown in Table 2, for our analysis for example by using the R language. For both SISR and OR-Tools we have that the associated $p$-value is lower than the significance level $\alpha$. We can thus reject the hypothesis $H_0$ that the average results of HGS-CVRP are similar to those of SISR and OR-Tools. Finally, by testing $H_1$ we again are able to reject the hypothesis, concluding that HGS-CVRP average results are statistically better than those obtained by SISR and OR-Tools on the $\mathbb{X}$ instances.
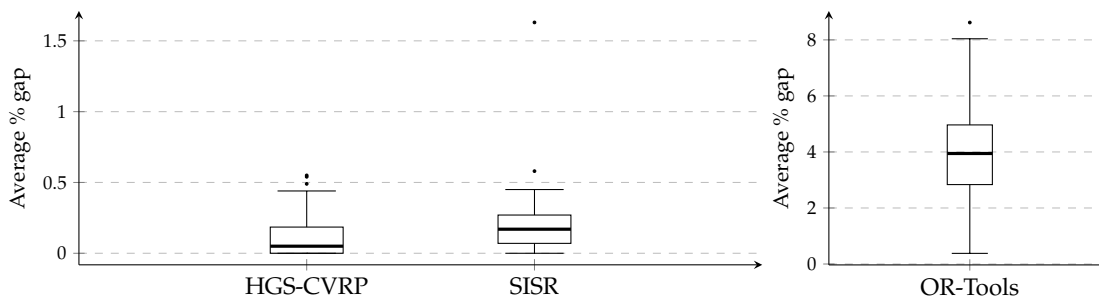


Figure 2: Average gaps obtained on the $\mathbb{X}$ instances by HGS-CVRP, SISR and OR-Tools. Note the different y-axis for OR-Tools. The thick line identifies the median value.

| | SISR | OR-Tools |
|---|---|---|
| $H_0$ | 8.27934e-06 | 3.95591e-18 |
| $H_1$ | 4.13967e-06 | 1.97796e-18 |

Table 2: $p$-values obtained when comparing HGS-CVRP with SISR and OR-Tools.

## 4  Conclusions

With this work we highlighted several challenges arising when comparing traditional and machine learning-based solution approaches for vehicle routing problems. Our aim was that of providing a set of reference guidelines that would help the machine learning community to produce computational studies that would be better appreciated by the operations research, and especially the vehicle routing, community. Finally, we conclude by referring the reader to the excellent overview on experimental analysis by Johnson (2002).

## References

[1]  Frank Wilcoxon. "Individual Comparisons by Ranking Methods". In: *Biometrics Bulletin* 1.6 (1945), pp. 80–83. ISSN: 00994987.

[2]  Olive Jean Dunn. "Multiple Comparisons Among Means". In: *Journal of the American Statistical Association* 56.293 (1961), pp. 52–64. ISSN: 01621459.

[3]  Marius M. Solomon. "Algorithms for the Vehicle Routing and Scheduling Problems with Time Window Constraints". In: *Operations Research* 35.2 (1987), pp. 254–265. eprint: `https://doi.org/10.1287/opre.35.2`

[4]  Hermann Gehring and Jörg Homberger. "A parallel hybrid evolutionary metaheuristic for the vehicle routing problem with time windows". In: *Proceedings of EUROGEN99*. Vol. 2. Citeseer. 1999, pp. 57–64.

[5]  David S. Johnson. "A theoretician's guide to the experimental analysis of algorithms". In: *Data Structures, Near Neighbor Searches, and Methodology: Fifth and Sixth DIMACS Implementation Challenges*. 2002, pp. 215–250.

[6]  Teodor Gabriel Crainic, Bertrand Le Cun, and Catherine Roucairol. "Parallel Branch-and-Bound Algorithms". In: *Parallel Combinatorial Optimization*. John Wiley & Sons, Ltd, 2006. Chap. 1, pp. 1–28. ISBN: 9780470053928. eprint: `https://onlinelibrary.wiley.com/doi/pdf/10.1002/9780470053928.ch1`.

[7]  Tobias Achterberg, Timo Berthold, and Gregor Hendel. "Rounding and Propagation Heuristics for Mixed Integer Programming". In: *Operations Research Proceedings 2011*. Ed. by Diethard Klatte, Hans-Jakob Lüthi, and Karl Schmedders. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 71–76. ISBN: 978-3-642-29210-1.

[8]  Paolo Toth and Daniele Vigo. *Vehicle Routing: problems, methods and applications*. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2014. eprint: `https://epubs.siam.org/doi/pdf/10.1137/1.97816`

[9]  Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. "Pointer Networks". In: *Advances in Neural Information Processing Systems*. Ed. by C. Cortes et al. Vol. 28. Curran Associates, Inc., 2015.

[10]  Hanjun Dai, Bo Dai, and Le Song. "Discriminative Embeddings of Latent Variable Models for Structured Data". In: *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*. ICML'16. New York, NY, USA: JMLR.org, 2016, pp. 2702–2711.

[11]    Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. http://www.deeplearningbook.org. MIT Press, 2016.

[12]    Irwan Bello et al. *Neural Combinatorial Optimization with Reinforcement Learning*. 2017. arXiv: 1611.09940 [cs.AI].

[13]    Marco Antonio Boschetti, Vittorio Maniezzo, and Francesco Strappaveccia. "Route relaxations on GPU for vehicle routing problems". In: *European Journal of Operational Research* 258.2 (2017), pp. 456–466. ISSN: 0377-2217.

[14]    Keld Helsgaun. *An Extension of the Lin-Kernighan-Helsgaun TSP Solver for Constrained Traveling Salesman and Vehicle Routing Problems: Technical report*. English. Roskilde Universitet, Dec. 2017.

[15]    Elias Khalil et al. "Learning Combinatorial Optimization Algorithms over Graphs". In: *Advances in Neural Information Processing Systems*. Ed. by I. Guyon et al. Vol. 30. Curran Associates, Inc., 2017.

[16]    Eduardo Uchoa et al. "New benchmark instances for the Capacitated Vehicle Routing Problem". In: *European Journal of Operational Research* 257.3 (2017), pp. 845–858. ISSN: 0377-2217.

[17]    Michel Deudon et al. "Learning Heuristics for the TSP by Policy Gradient". In: *Integration of Constraint Programming, Artificial Intelligence, and Operations Research*. Ed. by Willem-Jan van Hoeve. Cham: Springer International Publishing, 2018, pp. 170–181. ISBN: 978-3-319-93031-2.

[18]    Mohammadreza Nazari et al. *Reinforcement Learning for Solving the Vehicle Routing Problem*. 2018. arXiv: 1802.04240 [cs.AI].

[19]    Xinyun Chen and Yuandong Tian. "Learning to Perform Local Rewriting for Combinatorial Optimization". In: *Advances in Neural Information Processing Systems*. Ed. by H. Wallach et al. Vol. 32. Curran Associates, Inc., 2019.

[20]    Wouter Kool, Herke van Hoof, and Max Welling. "Attention, Learn to Solve Routing Problems!" In: *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

[21]    Laurent Perron and Vincent Furnon. *OR-Tools*. Version 7.2. Google, July 19, 2019.

[22]    Luca Accorsi and Daniele Vigo. *A Fast and Scalable Heuristic for the Solution of Large-Scale Capacitated Vehicle Routing Problems*. Tech. rep. University of Bologna, 2020.

[23]    Yoshua Bengio, Andrea Lodi, and Antoine Prouvost. *Machine Learning for Combinatorial Optimization: a Methodological Tour d'Horizon*. 2020. arXiv: 1811.06128 [cs.LG].

[24]    Francesco Cavaliere, Emilio Bendotti, and Matteo Fischetti. *An integrated local-search/set-partitioning heuristic for the Capacitated Vehicle Routing Problem*. Tech. rep. University of Padova, 2020.

[25]    Jan Christiaens and Greet Vanden Berghe. "Slack Induction by String Removals for Vehicle Routing Problems". In: *Transportation Science* 54.2 (2020), pp. 417–433. eprint: https://doi.org/10.1287/trsc.2019.0914.

[26]    Nina Mazyavkina et al. *Reinforcement Learning for Combinatorial Optimization: A Survey*. 2020. arXiv: 2003.03600 [cs.LG].

[27]    Artur Pessoa et al. "A generic exact solver for vehicle routing and related problems". In: *Mathematical Programming* 183.1 (Sept. 2020), pp. 483–523. ISSN: 1436-4646.

[28]    R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria, 2020.

[29]    Natalia Vesselinova et al. "Learning Combinatorial Optimization on Graphs: A Survey With Applications to Networking". In: *IEEE Access* 8 (2020), pp. 120388–120416. ISSN: 2169-3536.

[30]    Thibaut Vidal. *Hybrid Genetic Search for the CVRP: Open-Source Implementation and SWAP\* Neighborhood*. 2020. arXiv: 2012.10384 [cs.NE].

[31]  Yaoxin Wu et al. *Learning Improvement Heuristics for Solving Routing Problems*. 2020. arXiv: `1912.05784` `[cs.AI]`.

[32]  CVRPLIB. *Capacitated Vehicle Routing Problem Library*. visited on 2021-06-11. 2021.

[33]  André Hottung, Yeong-Dae Kwon, and Kevin Tierney. *Efficient Active Search for Combinatorial Optimization Problems*. 2021. arXiv: `2106.05126` `[cs.LG]`.

[34]  Wouter Kool, Herke van Hoof, Joaquim Gromicho, et al. *Deep Policy Dynamic Programming for Vehicle Routing Problems*. 2021. arXiv: `2102.11756` `[cs.LG]`.

[35]  OR-LIBRARY. *OR-Library is a collection of test data sets for a variety of Operations Research (OR) problems.* visited on 2021-07-03. 2021.

[36]  PassMark. *CPU Benchmarks*. visited on 2021-06-12. 2021.

[37]  TSPLIB. *TSPLIB is a library of sample instances for the TSP (and related problems) from various sources and of various types.* visited on 2021-07-03. 2021.

[38]  VRP-REP. *Collaborative open-data platform for sharing vehicle routing problem benchmark instances and solutions.* visited on 2021-06-11. 2021.