

Measurability and Reproducibility in University Timetabling Research: Discussion and Proposals

Andrea Schaerf and Luca Di Gaspero

Dipartimento di Ingegneria Elettrica, Gestionale e Meccanica
Università degli Studi di Udine
via delle Scienze 208, I-33100, Udine, Italy
{schaerf,l.digaspero}@uniud.it

Abstract. In this paper, we first discuss the level of compliance for timetabling research to two important research qualities, namely measurability and reproducibility, analyzing what we believe are the most important contributions in the literature. Secondly, we discuss some practices that, in our opinion, could contribute to the improvement on the two aforementioned qualities for future papers in timetabling research.

For the sake of brevity, we restrict our scope to university timetabling problems (exams, courses, or events), and thus we leave out other equally important timetabling problems, such as high-school, employee, and transportation timetabling.

1 Introduction

Thanks mainly to the PATAT conference series, researchers on timetabling problems have recently started to meet regularly to share experiences and results. This situation has the positive effect of generating both a common language and a common spirit that is the base ground for cross-fertilization of research groups in the timetabling community.

However, according to what we have seen at the recent PATAT conferences, the road for timetabling to become a well-established research community is still long. The main issue, in our opinion, is that most timetabling papers tend to describe the authors' specific problem and *ad hoc* solution algorithm without taking enough care of either the *measurability* or the *reproducibility* of the results. The reader is thus 'left alone' to judge the quality of the paper, and to understand what can be learned from it.

This issue is, to some extent, common to all the experimental areas of computer science and operations research, as clearly explained by Johnson in his seminal paper [17]. Nevertheless, we believe that this is particularly true in timetabling research, probably because of its shorter standing as a scientific community.

Regarding measurability (or comparability), we believe that several 'research infrastructures' are necessary in order to create the ground for truly measurable results. Specifically, they range from common formulations, to benchmark instances, to instance generators, to solution validators, and others. Related to it,

but somewhat complementary, is the issue of reproducibility. To this aim, beside the features just mentioned, it would be also necessary to create the conditions for sharing code and/or executables among researchers.

In this paper, we try to describe the main contributions with respect to these crucial qualities of experimental research in timetabling, and we also present some personal opinions on how to proceed to improve on them. For the sake of brevity, we restrict our scope to university timetabling problems (exams, courses, or events), and we leave out other equally important timetabling problems, such as high-school, employee and transportation timetabling. Nevertheless, to some extent, the proposed guidelines can have a broader application to all timetabling domains.

In detail, we first survey what, in our opinion, are the most important steps that have been pursued so far in timetabling research in terms of either measurability or reproducibility of results (Section 2). Secondly, we propose our personal ‘best practices’ for improving these two qualities in the timetabling research (Section 3). Our aim is to encourage both the authors to write research papers of high level in these important aspects and the reviewers to demand it when judging a paper.

2 Significant Contributions

In this section, we review the most significant contributions to the aim of creating the ground for the development of high quality measurable and reproducible research in timetabling. We first discuss the ‘standard’ problem formulations, the benchmark instances (datasets), and the related file formats adopted. Next, we move to the comparison methods proposed, such as competitions between algorithms and statistical tools. Finally, we discuss the issue of the objective validation of the proposed results.

2.1 Problem Formulations and Benchmark Instances

It is well known that timetabling problems vary not only from country to country, but also from university to university, and even in different departments of the same university the problem is not quite the same (see, e.g., [27]).

Nevertheless, throughout the years it has been possible to define common underlying formulations that could be used for the comparison of algorithms. In fact, a few basic formulations have become standards *de facto*, as they have been used by many researchers. Needless to say, standard formulations allow the researchers to compare their results and to co-operate for the solution. Furthermore, in some cases algorithms developed for more complex *ad hoc* formulations can be adapted to the basic standard ones so as to assess their objective quality.

For the Examination Timetabling problem (ETTP), Carter et al. [7] propose a set of formulations which differ from each other based on some components of the objective function. Carter also makes available a set of benchmark instances [6] extracted from real data, which represent a large variety of different situations. Formulations and benchmarks by Carter have stimulated a large body

of research, so that many researches (see, e.g., [4,8,15]) have adopted one of the formulations of Carter (or a variant of them, creating a new standard as well), tested on the benchmarks, and also added new instances. For more complex formulations, additional data have been added by other researchers, in an arbitrary but realistic way. At present, all available instances and the corresponding best results (only up to 2003, unfortunately) are published on the Web [20].

We call as Lecture Timetabling problem (LTTP), the problem of weekly scheduling a set of single lectures (or events). This problem differs from course timetabling (discussed below) because the latter is based on courses composed by multiple lectures, whereas lectures are independent. In fact, when a course is given in multiple lectures per week, some cost components are related to the way the lectures are placed in the week. In contrast, this concept is totally absent in LTTP. The LTTP differs also from ETTP because it has completely different objectives (e.g., no isolated event vs. spreading exams).

The LTTP has been discussed in [26] and it was the subject of the International Timetabling Competition ITC-2002 [23]¹. The formulation proposed for ITC2002 has also become quite standard, and many researchers have used it for their work (see, e.g., [9,19]). Twenty artificial instances were generated for the competition, and they are available from the ITC-2002 web page. In addition, a few others have been proposed (and made available via web) in [28].

Regarding the Course Timetabling problem (CTTP), which, as mentioned above, consists of the weekly scheduling of the lectures of a set of university courses, unfortunately no standard formulation has emerged from the community so far. To our knowledge, the only formulation available on the Web [14] together with a set of instances is the one proposed by ourselves in [16], along with four instances coming from the real cases (suitably simplified and made anonymous) in our university.

2.2 Data Format

For all the problems mentioned above, an important issue for the spreading of a formulation in the community is the data format. For all the formulations discussed above, the data format used is an ad hoc fixed-structure text-only one. For example, for ITC-2002 the input data comes in a single file containing the scalar values (number of events, rooms, room features, and students), followed by the elements of the input arrays, one per line. The output format follows the same idea. For the ETTP the input format is also rather ‘primitive’, with a fixed grammar and no formatting tags. Unfortunately, for this problem no output format has been specified in the original paper and the associated web site.

The use of fixed-structure formats makes it easier to parse the input from any computer language, and for any (naive) programmer, but may be more difficult

¹ Such a problem has often been referred to as CTTP (or UCTP), where C stands for course (and U for university); but we believe this is quite misleading, because it deals with isolated lectures/events, rather than courses composed by many lectures. Therefore we prefer for this problem the name LTTP.

to be maintained and checked. For example, it happened that Carter's ETPP instances were replicated incorrectly on other web sites. This was due to the presence of a few 'newline' characters added in the files, that led to different (less constrained) instances. As documented in detail in [25], this unfortunate episode has caused the publication of confusing results in some papers, and would have been avoided if a structured format had been used.

On the other hand, a structured format, such as XML, would be more suitable in terms of flexibility, extensibility, and maintenance, but it might hinder the work of researchers who cannot use it because of limited programming capabilities.

A few structured formats have been proposed in the literature, such as STTL [5,18] and TTML [22]. In [12], the authors go even beyond the language, proposing a multi-layer architecture for the specification and the management of timetabling problems. To our knowledge, however, these proposals have received limited attention so far in the academic community (although they are used in practical applications). This is probably due to the fact that researchers have normally little interest in the advantages of a structured language, and they prefer the quick-and-simple text-only version.

2.3 Comparison Methods and Competitions

The fair comparison of different algorithms and heuristics is well known to be a complex problem, and it has no simple and straightforward solution. In fact, in order to assess that an algorithm is 'better' than another it is necessary to specify not only the instances used, but also on which features they are compared under (e.g., quality of the objective function, success rate, speed, ...). The question gets even more complicated in presence of randomized/stochastic algorithms, which add a degree of non-determinism in the solution process.

For ITC-2002, the solution algorithms (provided as executables) were granted a maximum CPU time for their execution (based on a CPU benchmark, about 500 seconds on a recent computer) and they were evaluated only on the value of the objective function upon the 20 proposed instances. Unfeasible solutions were not considered, so that, in order to be admitted to the evaluation, participants had to find a feasible solution for all instances.

For stochastic algorithms, the participant had to ensure that their solver could produce the same solution when checked by the organization (by providing the seed of the random generator). In this situation, it is not clear how to apply the CPU time restriction and the choice of the organization was to grant the maximum time *for each single trial*. This was done to ensure reproducibility, although it had a drawback. The participants could take advantage of what we call the *Mongolian horde* approach: 'Run as many trials as you can and report only the best of all of them'. It is worth mentioning that in order to provide against the excessive use of the Mongolian horde approach, the competition organizers tested the best few algorithms also on unseen instances, and indeed the results were found to be broadly in line with the known instances.

To our knowledge, the ITC-2002 has been the sole attempt in this respect, and a new timetabling competition is scheduled for the second half of the 2007. All other comparisons are based on results published in the literature, which however often report only part of the necessary information (running times, number of trials, ...).

2.4 Result Validation

When some results are claimed in a research paper, the reader (or, more importantly, the reviewer) generally has to trust the author without any actual proof on the results. Although the possibility that the author is deliberately claiming a fake result is rare, cases in which the claimed results turned out to be wrong are relatively frequent. They are normally due to bugs in the code or misunderstandings in the formulation of the problem, typically the objective function.

For example, for the Graph Coloring problem, for the famous benchmark instance DSJC125.5 a 12-coloring solution has been claimed in 2002 (see [10]), whereas it has been successively proved that the minimum number of colors is 17.

Therefore the validation of the results claimed is clearly an important step toward the full reproducibility of the results. For the LTTP, in the ITC-2002, the validation of the results was done directly by the organizers, who asked all the participants to supply an executable that accepts a set of fixed command-line arguments.

For ETTP no validation tool has been available until very recently, and therefore validation has been based only on voluntary peer-to-peer interaction based on exchanges of solutions and values. Just before PATAT-2006, Rong Qu created a new web site [24] that allows the visitors to download an executable that validates ETTP solutions (using a raw fixed-structure output format). Up to now, the executable validates only solutions for the basic version of ETTP.

For our formulation of the CTTP, we have developed a web page [14] that allows other researchers to download the problem formulation, the data format, and the benchmark instances. More importantly, everybody is allowed also to upload and validate her/his own solutions, and to insert them among the results obtained for the specific instance. All results are automatically published on the web site along with the date and other information.

3 Proposals

In this section, we highlight some practices that, in our opinion, could contribute to the improvement on measurability and reproducibility for future papers in timetabling research. Part of what we propose here can be found also in [17], although we try to extract the advice of Johnson that we believe is best suited to the current state of timetabling research.

3.1 Statistically Principled Comparison

One of the key issues of performance measurement (often underestimated) concerns the methods to deal with the random nature of many techniques for obtaining a sound comparison of the different ones. In the practice, this issue is often neglected and just some tendency indicators of the stochastic variables, like mean values (and, more seldom, also standard deviations) in n runs (with $n \approx 10$), are provided. Furthermore, in a rather myopic view, these summary values are often advocated as the final word on the clear superiority of a technique over its competitors.

However, as is common knowledge in other research areas, when dealing with stochastic variables it is not correct to draw any conclusion only on the basis of single estimates, but a principled statistical analysis on the behavior of the algorithm is needed (see, e.g., [1,31]). Even in the simplest cases of comparison of two means, the analysis should include some kind of hypothesis testing (e.g., the t -test or the Mann–Whitney test for the parametric and the non-parametric case, respectively), that at least provides the reader with a probability measure of ‘confidence’ in the result. For more complex settings further analyses could be carried on and the statistical tool-case is plentiful of methods for correctly coping with several situations that arise in practice (see, e.g., [21]).

As an example, Birattari [2] has proposed a principled methodology for the comparison of stochastic optimization algorithms, called RACE, which comes out also as a software package for the R statistical software [3]. The RACE procedure, originally developed for the purpose of selecting the parameters of a single meta-heuristic, could be employed also in the case of the comparison of multiple algorithms by testing each of them on a set of trials. The algorithms that perform poorly are discarded and not tested anymore as soon as sufficient statistical evidence against them is collected. This way, only the statistically proven ‘good’ algorithms continue the race, and the overall number of tests needed to find the best one(s) is limited. Each trial is performed on the same randomly chosen problem instance for all the remaining configurations and a statistical test is used to assess which of them are discarded. The RACE procedure has been applied in the context of timetabling in [13].

It is worth noting that the statistical comparison of algorithms outlined in this section is based on the assumption of having full access to previous results (or, better, to the code) of the different techniques involved in the comparison. This is clearly related to the issue of reproducibility of results that, in our opinion, can be achieved by observing the guidelines described in the following.

3.2 Formulation, Data Format, Instances, and Results on the Web

As already mentioned, many papers in timetabling describe the modeling and the ad hoc solution of a new timetabling problem. For this kind of papers, in general we cannot expect that the authors make all the steps for obtaining full measurability and reproducibility such as, for example, publishing all the code. In fact, this would be quite a big job that would probably be too time-consuming for a researcher, beside possible employer’s concerns. Nevertheless, we believe

that there are a few actions that could contribute in these respects, which are not too expensive in terms of human work.

First, the authors must state the problem clearly and exhaustively. If this is not possible in the paper for space reasons, the full formulation should be posted in an accompanying web site. Secondly, the authors should also post in the web site all the instances considered in the study (hiding identities for privacy reasons, if necessary), along with all the necessary information accompanying them: data format, algorithms, results, and running times. Finally, the authors should post also the files containing their best solutions, so that other researcher can verify the actual results, and possibly use that solutions for further studies and improvements.

These actions would ensure comparability with the results on future research by other researchers or also by the same authors.²

3.3 Web-Based Problem Management System

Nowadays it is very common to see web sites that describe all aspects of either a specific problem, see e.g. [11,29], or a research area [30]. These web sites normally exhibit references to papers, people, problem formulations, and benchmark instances, and supply other information.

Web sites are surely very useful for the community, and their presence is crucial for the quality of the research. Nevertheless, we believe that there is a further step to be made to this regard. Inspired by the well-known concept of CMS (*content management system*), we envision the idea of developing what we would call PMS (*problem management system*). A PMS is a web application (rather than a web site) that should allow the users to interact with the application performing all the following tasks:

Add results: New results are first validated, and then possibly inserted in the database along with the time-stamp and other user-supplied information.

Add instances: Instances can be inserted at any moment. Researchers that are interested in the problem can be automatically informed by email of this kind of event.

Manage instance generation: Newly generated instances can be created automatically by users through interaction with an instance generator.

Analyze instances and results: Instances and results can be analyzed automatically so as to produce important indicators: constrainedness, similarity to other instances or other results, etc.

Add general information: People, references, links, code, and other information can be added. Links would be validated periodically in an automatic way, and broken ones can be removed. References can also be imported from other sites.

Translate data: Input and output data can be translated so that coherent data can be proposed in different formats to the community (including both fixed-structure and XML-based ones).

² Many researchers – including ourselves! – have experienced the frustration of losing their solutions (or other data) for some of the problems they have worked on.

Organize on-line competitions: Competitions on specific instances and with registered participants and fixed deadlines can be organized semi-automatically. Results can be reported immediately.

Visualize: Solutions can be visualized in graphical form to give an immediate picture of the features and the violations.

Maintain a discussion forum: A simple discussion forum about the problem can be maintained along with the site. Messages would be organized and displayed as in usual on-line forums (threads, date, ...).

The interesting point is that information posted through the PMS would get on-line immediately in an automatic way. Obviously, a PMS needs to provide against possible abuses and malicious behavior, and therefore some of the actions mentioned above would need the approval of the administrator before becoming effective. For most operations, this however would be just a Yes/No button, so that the administrator will hopefully operate in short time.

The PMS would also maintain historical data (through versioning systems), in such a way to be able to retrieve information eliminated by updates and deletions.

Acknowledgements. We wish to thank Marco Chiarandini for fruitful discussions about measurability and reproducibility of research results. We also thank the anonymous referees for their comments that helped us to improve the paper.

References

1. Barr, R., Golden, B., Kelly, J., Resende, M., Stewart, W.: Designing and reporting on computational experiments with heuristic methods. *Journal of Heuristics* 1, 9–32 (1995)
2. Birattari, M.: The problem of tuning metaheuristics as seen from a machine learning perspective. Ph.D. Thesis, Université Libre de Bruxelles, Belgium (2004)
3. Birattari, M.: The RACE Package (April 2005)
4. Burke, E.K., Newall, J.: A multi-stage evolutionary algorithm for the timetable problem. *IEEE Transactions on Evolutionary Computation* 3, 63–74 (1999)
5. Burke, E.K., Pepper, P., Kingston, J.: A standard data format for timetabling instances. In: Burke, E.K., Carter, M. (eds.) PATAT 1997. LNCS, vol. 1408, pp. 213–222. Springer, Heidelberg (1998)
6. Carter, M.W.: Carter’s test data (2005) (viewed March 13, 2007) (updated June 7, 2005), <ftp://ftp.mie.utoronto.ca/pub/carter/testprob/>
7. Carter, M.W., Laporte, G., Lee, S.Y.: Examination timetabling: algorithmic strategies and applications. *Journal of the Operational Research Society* 74, 373–383 (1996)
8. Casey, S., Thompson, J.: Grasping the examination scheduling problem. In: Burke, E.K., De Causmaecker, P. (eds.) PATAT 2002. LNCS, vol. 2740, pp. 232–244. Springer, Heidelberg (2003)
9. Chiarandini, M., Birattari, M., Socha, K., Rossi-Doria, O.: An effective hybrid approach for the university course timetabling problem. *Journal of Scheduling* 9, 403–432 (2006)

10. The website of graph coloring and its generalization (2004) (viewed March 13, 2007), <http://mat.gsia.cmu.edu/COLOR04>
11. Culberson, J.: Graph coloring page (2006) (viewed March 13, 2007) (updated March 31, 2004), <http://www.cs.ualberta.ca/~joe/Coloring/>
12. De Causmaecker, P., Demeester, P., Lu, Y., Vanden Berghe, G.: Using Web standards for timetabling. In: Burke, E.K., De Causmaecker, P. (eds.) PATAT 2002. LNCS, vol. 2740, pp. 238–257. Springer, Heidelberg (2003)
13. Di Gaspero, L., Chiarandini, M., Schaerf, A.: A study on the short-term prohibition mechanisms in tabu search. In: ECAI-2006. Proceedings of the 17th European Conference on Artificial Intelligence, Riva del Garda, Italy, pp. 83–87.
14. Di Gaspero, L., Fontanel, A., Schaerf, A.: Educational Timetabling @UniUd (2006) (viewed March 13, 2007) (updated May 30, 2006), <http://www.diegm.uniud.it/satt/projects/EduTT/>
15. Di Gaspero, L., Schaerf, A.: Tabu search techniques for examination timetabling. In: Burke, E., Erben, W. (eds.) PATAT 2000. LNCS, vol. 2079, pp. 104–117. Springer, Heidelberg (2001)
16. Di Gaspero, L., Schaerf, A.: Multi-neighbourhood local search with application to course timetabling. In: Burke, E.K., De Causmaecker, P. (eds.) PATAT 2002. LNCS, vol. 2740, pp. 262–275. Springer, Heidelberg (2003)
17. Johnson, D.S.: A theoretician’s guide to the experimental analysis of algorithms. In: Goldwasser, M.H., Johnson, D.S., McGeoch, C.C. (eds.) Data Structures, Near Neighbor Searches, and Methodology: 5th and 6th DIMACS Implementation Challenges, pp. 215–250. American Mathematical Society, Providence, RI (2002), available from <http://www.research.att.com/~dsj/papers.html>
18. Kingston, J.H.: Modelling timetabling problems with STTL. In: Burke, E., Erben, W. (eds.) PATAT 2000. LNCS, vol. 2079, pp. 309–321. Springer, Heidelberg (2001)
19. Kostuch, P.: The university course timetabling problem with a three-phase approach. In: Burke, E.K., Trick, M.A. (eds.) PATAT 2004. LNCS, vol. 3616, pp. 109–125. Springer, Heidelberg (2005)
20. Merlot, L.: Public Exam Timetabling Data Sets (2005) (viewed March 13, 2007) (updated October 13, 2003), <http://www.or.ms.unimelb.edu.au/timetabling>
21. Montgomery, D.C.: Design and Analysis of Experiments, 6th edn. Wiley, New York (2005)
22. Özcan, E.: Towards an XML-based standard for timetabling problems: TTML. In: Kendall, G., Burke, E.K., Petrovic, S., Gendreau, M. (eds.) MISTA 2003. Proceedings of the 1st Multidisciplinary International Conference on Scheduling: Theory and Applications, pp. 163–185. Springer, Berlin (2005)
23. Paechter, B., Gambardella, L.M., Rossi-Doria, O.: International Timetabling Competition Webpage (2003) (viewed March 13, 2007) (updated July 10, 2003), <http://www.idsia.ch/Files/ttcomp2002/>
24. Qu, R.: The exam timetabling site (2006) (viewed March 13, 2007) (updated July 8, 2006), <http://www.cs.nott.ac.uk/~rxq/ETTP.htm>
25. Qu, R., Burke, E.K., McCollum, B., Merlot, L., Lee, S.Y.: The state of the art of examination timetabling. Technical Report NOTTCS-TR-2006-4, School of CSiT, University of Nottingham, UK (2006)
26. Rossi-Doria, O., Sampels, M., Birattari, M., Chiarandini, M., Dorigo, M., Gambardello, L.M., Knowles, J., Manfrin, M., Mastrolilli, M., Paechter, B., Paquette, L., Stützle, T.: A comparison of the performance of different metaheuristics on the timetabling problem. In: Burke, E.K., De Causmaecker, P. (eds.) PATAT 2002. LNCS, vol. 2740, pp. 329–351. Springer, Heidelberg (2003)

27. Schaerf, A.: A survey of automated timetabling. *Artificial Intelligence Review* 13, 87–127 (1999)
28. Socha, K., Knowles, J., Sampels, M.: A MAX-MIN Ant System for the University Timetabling Problem. In: Fieker, C., Kohel, D.R. (eds.) *Algorithmic Number Theory*. LNCS, vol. 2369, Springer, Heidelberg (2002), Data available from: <http://iridia.ulb.ac.be/~msampels/tt.data/>
29. Trick, M.: Challenge Traveling Tournament Instances, web page (2005) (viewed March 13, 2007) (updated October 22, 2006), <http://mat.gsia.cmu.edu/TOURN/>
30. The web site of the EURO working group on automated timetabling (WATT) (2002) (viewed March 13, 2007) (updated February 21, 2007), <http://www.asap.cs.nott.ac.uk/watt/>
31. Xu, J., Chiu, S., Glover, F.: Fine-tuning a tabu search algorithm with statistical tests. *International Transactions on Operational Research* 5, 233–244 (1998)