

Constraints

Optimal Methods for Resource Allocation and Scheduling - A Cross-Disciplinary Survey

--Manuscript Draft--

Manuscript Number:	CONS213R1
Full Title:	Optimal Methods for Resource Allocation and Scheduling - A Cross-Disciplinary Survey
Article Type:	Survey Paper (Thomas Schiex)
Keywords:	Scheduling; Resource Allocation; Constraint Programming; Operations Research; Hybrid Algorithms
Abstract:	<p>Classical scheduling formulations typically assume static resource requirements and focus on deciding when to start the problem activities, so as to optimize some performance metric. In many practical cases, however, the decision maker has the ability to choose the resource assignment as well as the starting times: this is a far-from-trivial task, with deep implications on the quality of the final schedule.</p> <p>Joint resource assignment and scheduling problems are incredibly challenging from a computational perspective. They have been subject of active research in Constraint Programming (CP) and in Operations Research (OR) for a few decades, with quite difference techniques. Both the approaches report individual successes, but they overall perform equally well or (from a different perspective) equally bad. In particular, despite the well known effectiveness of global constraints for scheduling, comparable results for joint filtering of assignment and scheduling variables have not yet been achieved. Recently, hybrid approaches have been applied to this class of problems: most of them work by splitting the overall problem into an assignment and a scheduling subpart; those are solved in an iterative and interactive fashion with a mix of CP and OR techniques, often reporting impressive speed ups compared to pure CP and OR methods.</p> <p>Motivated by the success of hybrid algorithms on resource assignment and scheduling, we provide a cross-disciplinary survey on the topic, including CP, OR and hybrid approaches. The main effort is to identify key modeling and solution techniques, so that they may be applied in the construction of new hybrid algorithms, or they may provide ideas for devising novel filtering methods (possibly based on decomposition, or on alternative representations of the domain store). In detail, we take a constraint-based perspective and, according to the equation $CP = \text{model} + \text{propagation} + \text{search}$, we give an overview of state-of-art models, propagation/bounding techniques and search strategies.</p>

Optimal Methods for Resource Allocation and Scheduling: a Cross-disciplinary Survey

– Accompanying Letter –

Michele Lombardi Michela Milano

September 26, 2011

Dear editor,

this is an accompanying letter to the revised version of our manuscript “Optimal Methods for Resource Allocation and Scheduling: a Cross-disciplinary Survey”. We bring our apologies for the delay, but we decided to address thoroughly the received comments, despite just a minor revision was required. We thank the reviewers for pointing out missing references and for having highlighted the main weaknesses of the first version of this paper, pushing us to heavily revise it and to improve its quality. In particular:

- We clarified the main purpose of the survey: in particular, we considerably modified the abstract, the introduction and the conclusion section.
- The material from former Section 2 has been split into distinct sections. In particular, Section 2 in the current paper is entirely dedicated to the problem definition, while modeling approaches are described in Section 3. As a side-effect, some discussion about the impact of different problem variants on complexity has been moved to Section 4 and following ones. Most notably, this was done for “regular cost functions” and “generalized precedence constraints”.
- We revised Section 3, improving in particular the discussion on CP modeling techniques and hybrid/decomposition-based approaches.
- We re-organized Section 4 (formerly Section 3): CP filtering techniques are now classified as temporal, logical and resource reasoning (depending on the aspect they mainly focus on). Moreover, we added several of the references mentioned by the reviewers.
- We simplified and slightly re-organized Section 5 (formerly Section 4), about bounding and dominance rules. In particular, the discussion on lower bounds (Section 5.4, formerly 4.1) is now much more focused on allocation and scheduling problems.

- Section 6.1 (formerly 5.1) about search techniques in CP has been restructured and strategies with strong similarities are now discussed together.
- The former Section 5.3 (on Logic based Benders' Decomposition) has been extended, in part by embedding some of the discussion from former Section 2.4 and in part by adding new material. The result is the novel Section 7.
- We revised the whole text, trying to improve fluency and quality of the presentation.
- For each technique described in the paper, we tried to point out (as much as possible) strengths, weaknesses and connections to other methods.
- We removed or reduced the descriptions of secondary techniques, leaving more room for the discussion of the most important ones.
- We cleaned-up the references, eliminating redundant and non-sufficiently-related papers.

Additionally, we have corrected all reported typos. The resulting paper exceeds the typical length limit (31 pages rather than 30): given the amount of references, we hope this is still acceptable. Nevertheless, we are ready to shorten it to 30 pages, in case it is requested.

Best regards,

Michele Lombardi, Michela Milano

Response to Reviewer 1

In the description of the different models, I think there is a wrong claim about the model with optional activities and time intervals: “Note however this technique requires the introduction of an exponential number of variables in case of independent resource groups”. In this type of model, the same interval variable “a” can be the master of several alternative constraints. If to take the example of the workers and factory rooms from Figure 1, one could model an activity a that requires a worker of type r0 or r1 and 1 factory room r2 or r3 as:

```
alternative(a, [ar0, ar1]);
alternative(a, [ar2, ar3]);
```

where a_{ri} are optional interval variables representing the execution of a on resource r_i.

This model does not require an exponential number of intervals. If needed, compatibility constraints can then be modeled by logical constraints (example: presenceOf(a_{r0}) \Rightarrow presenceOf(a_{r2})).

Of course, another way to model is to create one optional interval for each element of the cartesian product of the different types of resources. This model

is indeed exponential and it is interesting only if the number of possible configurations is not too large, typically when resources are very dependent from each other (this relates with the notion of activity “mode”).

An interest of optional interval variables is to allow for both types of models. I think this part of the paper should be changed, it is mentioned twice on pages 6 and 7.

We agree with the reviewer, despite the claim actually holds if Generalized Precedence Constraints are not supported (in particular, maximal time lags are needed to enforce synchronization of start/end times). This is now clearly discussed in Section 3.1.2 (when talking about alternative activities and interval variables) and in Section 3.2 (about the MRCPSp).

It would be interesting to consider some of the following studies in the review. These papers deal with models involving scheduling and resource allocation.

Separate comments are given for each proposed paper.

(1) *Scheduling of landing of aircrafts on several runways (includes runway allocation)*

1. A. T. Ernst, M. Krishnamoorthy, and R. H. Storer. *Heuristic and exact algorithms for scheduling aircraft landings.* Networks, 34(3):229-241, 1999.
2. A. T. Ernst, M. Krishnamoorthy. *Algorithms for Scheduling Aircraft Landings.* 2001.

We added reference 1, which is now discussed in several parts of the paper.

(2) *Resource allocation in the context of audit scheduling problems*

1. Peter Brucker, Doris Schumacher. *A new tabu search procedure for an audit-scheduling problem.* Journal of Scheduling, 2: 157-173, 1999.

Omitted, since the survey focus is on tree-search based approaches.

(3) *Flexible jobshop (machine allocation)*

1. Mastrolilli M., Gambardella L.M. , (1998) *Effective Neighborhood Functions for the Flexible Job Shop Problem,* Journal of Scheduling, Volume 3, Issue 1, 2000. Pages: 3-20

Briefly mentioned in Section 1, since it is not based on tree-search approach.

(4) *Multi Skill Project Scheduling Problem (manpower allocation)*

1. C. Pesson, O. Bellenguez-Morineau, E. Neron. *Multi-skill Project Scheduling Problem and Total Productive Maintenance.* Proceedings MISTA 2007.

This is now cited and discussed in the paper.

(5) *Oversubscribed scheduling problem with resource allocation*

1. L. A. Kramer and L. V. Barbulescu and S. F. Smith. *Understanding Performance Tradeoffs in Algorithms for Solving Oversubscribed Scheduling*. Proceedings AAAI 2007.

Omitted, since the paper main focus is on evaluation rather than on solution techniques.

(6) *Vehicle allocation in Routing problems*

Vehicle routing problems often involve a vehicle allocation component. And some of these problems are at the frontier between routing and scheduling. For instance Capacitated Vehicle Routing Problem with Time Windows.

Including Vehicle Routing Problems would make the survey way too large; therefore we just added a reference to a Vehicle Routing book in the introduction.

p4: “*the most common problem objective is the makespan*”. Well, that is true for academical problems but in industrial applications tardiness and allocations costs often play a more important role.

We agree: this is now stated clearly in Section 2.1.3.

p10: In “*Other Objective Functions*” you mention weighted tardiness. May be this type of cost should better fit in section 2.5.1 in “*Time Based Objectives*” ?

The whole section has been restructured: we hope it is now more clear and logically organized.

p13: section 3.2. The last sentence is surprising. Usually one uses a cumulative model (sometimes as a redundant constraint on top of the alternatives) precisely because it allows a stronger propagation while actual resource assignment is still not performed. The sentence seems to suggest the opposite.

True. In detail, the performed filtering is complementary: redundant resources tend to be useful as long as many activities are still unassigned. This is now stated clearly in Section 4.3.

In the references, I think the date for paper [76] (Kwok) is 1996 rather than 2002.

Fixed: it was just a typo.

Response to Reviewer 2

Good surveys are very hard to write as they very easily degenerate into an annotated bibliography or, even worse, a list of papers. The authors have attempted a categorization which is both necessary and reasonably successful. Unfortunately, within each of these categories (e.g. propagation, search, etc.) the survey tends to resemble more of an annotated bibliography. For example, in the search section there seem to be possible connections amongst the topics discussed: the resource envelopes of the precedence constraint posting seem to have a similar flavour to the texture measurement of the heuristic commitment techniques; left-justified random appears to be similar to schedule or postpone. My point isn't that these methods are identical but that there are underlying ideas that have been developed independently and that one very valuable use of a survey is to distill these ideas. Unfortunately, such a distillation is not done.

As mentioned in the introduction of this accompanying letter, we revised the whole paper, trying to point out the best performing techniques and to outline the connections between similar ideas. For sure there is still much to dig up, but we hope our effort goes in the direction suggested by your comment.

Secondly, a main part of the value of a survey is the isolation and highlighting of the really significant papers. What is included is as important as what is left out and it is the judgement of the authors on this issue which provides value. Again, unfortunately, I did not get much of a sense about which papers in each theme were seminal and/or important as opposed to making incremental contributions.

We modified the paper so as to explicitly identify papers where a new technique has been introduced and papers where some existing technique is just applied or extended. We tried hard to reduce the number of redundant references, keeping a few of them for each topic: as a result, despite the addition of the new papers suggested by the reviewers, the bibliography dropped from 127 to 91 entries.

Many of the different papers reviewed solve different problems. On one hand many of the CP-based approaches solve simpler problems (without allocation) [for example most of the search techniques] while many of the OR approaches solve more general problems (i.e. the multi-mode scheduling). This dichotomy ends up confusing the reader as it is never clear if (or how) a given technique (such as precedence constraint posting, to pick one) has been applied to problems with resource allocation. Conversely, it is sometimes unclear what the more general problems that the OR techniques are solving actually are and how such techniques would do if applied to basic allocation/scheduling problems. Not only is this organization somewhat confusing, it also is misleading in that it seems to indicate that many more pieces of work have focused on allocation/scheduling problems (i.e., in CP scheduling) than is really the case. I suggest that if the focus of the survey is truly combined allocation/scheduling problems, much of the work that does not explicitly attack such problems should be discarded. It may

*be that a short section on "foundational" techniques (e.g., the standard non-optimal propagators) could be provided but I think much of the detail can be skipped as it is not directly relevant (except in a trivial way: assign all resources then do normal scheduling). Similarly, the paper could have a shorter section on techniques that have been applied to *more* general problems but without some indication of how these techniques work on allocation/scheduling problems, they should be de-emphasized.*

We drastically reduced the number of citations not directly tackling joint resource assignment and scheduling problems. We did not introduce a section for foundational techniques, as it would make the discussion too fragmented. Conversely, we modified the paper so as to identify clearly all references on pure scheduling.

As for MRCPS approaches, we think they are worth the attention they receive in the paper, for a number of reasons: (1) the main purpose of this paper is to survey useful techniques for solving assignment and scheduling problems, and techniques developed for the MRCPS indeed are useful for that. (2) Since this paper is submitted to "Constraints" and puts some emphasis on the effectiveness of hybrid algorithms, exposing methods from the OR community besides CP ones provides a very valuable contribution. (3) Recent CP techniques (i.e. alternative activities and conditional time interval variables) are able to model multiple modes and even more general problems. Since we moved the discussion on problem variants in Section 2, immediately after the problem definition, it should now be clear that the scope of the paper is sufficiently broad to include the MRCPS.

One argument for this paper appearing in Constraints is that scheduling is generally seen (within the CP community) as a success story. The paper points out significant other work that CP schedulers may not be aware of as well as addressing problem that CP doesn't do a particularly good job on. On the other hand, the core theme of the paper is a type of scheduling problem that probably will not interest a large portion of the Constraints community. I would expect to see this paper in something like the Journal of Scheduling. I do not want to over-state my opinion here as I believe that Constraints certainly should be broad enough for this type of paper. However, the authors might be much more mileage from the paper if it appeared in a forum where more people directly interested in scheduling would read it. Furthermore, if the authors really want to aim at Constraints, I believe that they could have done more to make the paper more directly relevant. It would be nice to see some motivation at the beginning of the paper answering the question: why should one who is not interested in scheduling read this paper? I think there are a number of arguments that could be made here (i.e, the "success" story in scheduling but the lack of success in resource allocation/scheduling serving as a challenge to the CP community).

We agree: in particular, in our opinion resource allocation and scheduling problems provide a nice benchmark for devising hybrid algorithms (they have already pushed the development of many of them) and experimenting novel

filtering techniques (e.g. based on decomposition, SAT hybridization or alternative representations of the domain store). This is now stated in the abstract, in the introduction, in the conclusion and at the beginning of Section 4.

The English level is not very high. While the paper is generally understandable, there is hardly a paragraph which doesn't contain English mistakes. The first line of the abstract is included here: "Resource Allocation and Scheduling problems consist into assigning over time resources from a candidate pool to a set of activities;" → "consists of" and the second part of the sentence "assigning over ..." is basically incomprehensible. Again, I do not think this should be overstated as writing academic papers in a second language is not easy. However, the level here is below what I would expect as a submission to Constraints especially as most of the problems can be fixed by a native-English editor. There are even typos that should have been caught by a spell-checker ("aplped").

We revised the whole paper, trying to improve the language fluency and correctness.

The paper does a good job of covering relevant papers though doesn't seem to express judgement in pointing out really crucial papers and de-emphasizing or omitting those that make less of a contribution. The only major gap that I saw was the work in COSYTEC on the `diffn` constraint. Indeed, there is a substantial amount of propagation based work (much of it due to Beldiceanu) on rectangle packing that is a generalization of resource allocation/scheduling problems. Unfortunately, some of it appears in obscure workshops. Perhaps, contacting some of the responsible authors directly (Beldiceanu, Simonis, etc.) will help to find these works.

As previously mentioned, we cleaned-up the bibliography by removing minor, non sufficiently related and redundant papers. We have included new references on filtering algorithms for cumulative constraints in the presence of resource assignment choices.

We did not include a discussion of the `diffn` constraint, since cumulative scheduling is a *relaxation* of rectangular packing (see e.g. [1]): some instances may admit a feasible schedule even if the corresponding packing problem has no solution, so that using packing constraint to formulate a scheduling problem incurs the risk of loosing optimality. This is now mentioned in the paper in Section 3.2.1, when discussing a variant of the disjunctive MILP model.

As a final general comment: there are a number of notions (e.g., "Resource Strength", "Force Directed Scheduling") that are mentioned without definition. I do not think doing this really contributes to a survey - if the ideas are important enough to be mentioned, they should be defined.

Non-defined terms in the paper have been either removed, or given a definition or described in an intuitive fashion.

Response to Reviewer 3

The coverage of both the CP and OR literatures is quite thorough. This should be a valuable resource for the CP audience to which it is directed, particularly the material on OR methods, which may be unfamiliar. The CP material would be of equal interest to the OR community, and it is a pity that the paper will not reach this audience. I suggest that the authors consider writing a version of the paper for an OR journal, perhaps with some more tutorial material on CP methods. CP-based and hybrid scheduling methods are largely unknown to the OR world. One possible outlet would be INFORMS Journal on Computing.

We will certainly consider your suggestion.

I collected some citations on decomposition methods that may be useful: [2, 3, 4, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 19, 22, 23, 24, 25, 26]. Logic-based Benders was introduced in [15]. You may want to mention CP-based branch and price [8, 17, 27, 28]. More references appear in [8]. Also some references on large neighborhood search, briefly described at the end of the paper. The idea is introduced in [20] and surveyed in [18]. There is also very large neighborhood search [1], although I'm not sure how it differs. Hybrids of local search and CP are surveyed by [21] and hybrid metaheuristics by [5].

In order to keep the survey more focused, as requested by other reviewers, we checked all the suggested references and added the most related ones. The paper still contains no-mention to CP branch-and-price since, to the authors knowledge, has never been applied to resource allocation and scheduling. Heuristic methods are unfortunately outside of the scope of this paper: due to their importance, however, they are mentioned in the introduction, where the reader is invited to check some key references to get more details.

References

- [1] Nicolas Beldiceanu, Mats Carlsson, Sophie Demassey, and Emmanuel Poder. New filtering for the cumulative constraint in the context of non-overlapping rectangles. *Annals of Operations Research*, 184(1):27–50, March 2011.

Noname manuscript No.
(will be inserted by the editor)

1
2
3
4 **Optimal Methods for Resource Allocation and Scheduling:**
5 **a Cross-Disciplinary Survey**
6
7
8 **Michele Lombardi · Michela Milano**
9
10
11
12
13
14
15 Submitted: December 2010
16
17

18 **Abstract** Classical scheduling formulations typically assume static resource requirements and focus on deciding when to start the problem activities, so as to optimize some performance metric. In many practical cases, however, the decision maker has the ability to choose the resource assignment as well as the starting times: this is a far-from-trivial task, with deep implications on the quality of the final schedule.

19 Joint resource assignment and scheduling problems are incredibly challenging from
20 a computational perspective. They have been subject of active research in Constraint
21 Programming (CP) and in Operations Research (OR) for a few decades, with quite
22 difference techniques. Both the approaches report individual successes, but they over-
23 all perform equally well or (from a different perspective) equally bad. In particular,
24 despite the well known effectiveness of global constraints for scheduling, comparable
25 results for joint filtering of assignment and scheduling variables have not yet been
26 achieved. Recently, hybrid approaches have been applied to this class of problems:
27 most of them work by splitting the overall problem into an assignment and a schedul-
28 ing subpart; those are solved in an iterative and interactive fashion with a mix of CP
29 and OR techniques, often reporting impressive speed ups compared to pure CP and
30 OR methods.
31
32

33 Motivated by the success of hybrid algorithms on resource assignment and schedul-
34 ing, we provide a cross-disciplinary survey on the topic, including CP, OR and hybrid
35 approaches. The main effort is to identify key modeling and solution techniques, so that
36 they may be applied in the construction of new hybrid algorithms, or they may pro-
37 vide ideas for devising novel filtering methods (possibly based on decomposition, or on
38 alternative representations of the domain store). In detail, we take a constraint-based
39 perspective and, according to the equation $CP = \text{model} + \text{propagation} + \text{search}$, we
40 give an overview of state-of-art models, propagation/bounding techniques and search
41 strategies.
42
43

44 **Keywords** Scheduling, Resource Allocation, Constraint Programming, Operations
45 Research, Hybrid Algorithms
46
47

48 M. Lombardi, M. Milano
49 DEIS, University of Bologna, Viale del Risorgimento 2, 40136 Bologna E-mail:
50 {michele.lombardi2,michela.milano}@unibo.it
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Resource Allocation and Scheduling

Following a widely accepted definition, scheduling involves “allocating scarce resources to activities over time”. Classical scheduling formulations¹ typically assume static resource requirements and focus on deciding when to start the problem activities, so as to optimize some performance metric. Sometimes, however, the decision maker has the ability to choose the resources for the execution of each problem activity; this is a far-from-trivial task, with deep implications on the quality of the final schedule. Despite considerably less popular than classical scheduling in the research literature, joint resource assignment and scheduling problems are very frequent in practice: according to [50,11], many real world settings feature *optional activities* which can be left non-executed (usually with an impact on costs) or *alternative recipes* to execute an activity, with each recipe corresponding to a possible assignment of execution resources or to a different set of sub-activities.

Here, we take into account a variety of scheduling problems, having as a common feature user-decided resource requirements. Reported applications include batching and scheduling in chemical plants with several reactors [83]; scheduling aircraft landings to multiple runways [10]; simultaneous scheduling of maintenance activities requiring different skills [14]; compilation of computer programs on heterogeneous Very Large Instruction Word (VLIW) architectures [60]. Problems of this class do also arise in the optimization of parallel applications on multi-core platforms [52].

Scheduling problems are well-known to be NP-hard and computationally challenging: not surprisingly, introducing resource assignment decisions drastically increases the complexity. As an example, the well-known PSPLIB benchmarks [49] include classical scheduling problems (Resource Constrained Project Scheduling Problems – RCPSP) as well as problems involving resource assignment decisions (Multi-mode Resource Constrained Project Scheduling Problems – MRCPS): while RCPSP instances with up to 120 activities are solved by powerful hybrid CP/SAT methods [74], finding optimal solutions to 30 activity instances of the MRCPS is still considered a very challenging task [91].

Constraint Programming (CP) can claim a success story with pure scheduling [3], providing an elegant and versatile modeling framework, based on effective and efficient filtering algorithms and search strategies. Significant effort has been put into extending CP models and solution techniques for allocation and scheduling problems; unfortunately the obtained results are not comparable to those of pure scheduling, the main reasons are the very large search space and lack of effective techniques for joint filtering of assignment and scheduling variables.

Operations Research (OR) scientists have also considered resource allocation and scheduling problems (see the survey by Brucker et al. [18]), mainly under the flag of the MRCPS, introduced in the late 70’s [31]. Taking advantage of a strong mathematical basis, they have identified relevant problem properties and powerful optimization-based filtering rules (dominance rules). Nevertheless, practical experience shows no sharp dominance of OR techniques over CP ones: both approaches report individual successes, but they overall perform equally well or (from a different perspective) equally bad.

Starting from 2001, hybrid CP/OR approaches have been applied to resource allocation and scheduling problems: to the best of the authors’ knowledge, works [47, 82, 40] are the earliest, most relevant examples; the mentioned approaches are mainly

¹ E.g. the Resource Constrained Project Scheduling Problem or Job Shop Scheduling.

1 based on Logic Based Benders' Decomposition (introduced in [45] and formalized in
 2 [42]) and report orders-of-magnitude speed-ups compared to pure CP and pure OR
 3 methods. Intuitively, since allocation and scheduling problems result from the com-
 4 position of an assignment and a scheduling component, hybrid algorithms have the
 5 opportunity to use the most effective (heterogeneous) technique to target each prob-
 6 lem part. (e.g. Mixed Integer Linear Programming – MILP – for the assignment and
 7 CP for scheduling).

8 Motivated by the success of hybrid methods on resource assignment and scheduling,
 9 we provide a cross-disciplinary survey on the topic, including CP, OR and hybrid
 10 approaches. *The main effort is to identify key modeling and solution techniques, so*
 11 *that they may be applied in the construction of new hybrid algorithms, or they may*
 12 *provide ideas for devising novel filtering methods* (e.g. based on decomposition or on
 13 alternative representations of the domain store). In detail, we take a constraint-based
 14 perspective and, according to the equation $CP = \text{model} + \text{propagation} + \text{search}$, we
 15 give an overview of state-of-art models, propagation/bounding algorithms and search
 16 strategies.

17 In particular, we focus on approaches making some use of tree search, as the tech-
 18 niques they apply are more easily portable in a CP context. This choice excludes local
 19 search, greedy heuristics and all meta-heuristic methods; due to the problem complex-
 20 ity, those are however very important approaches: work [85] reports a comprehensive
 21 list of previous heuristic and meta-heuristic approaches for the MRCSP and describes
 22 a novel genetic algorithm; [35] describes a metaheuristic based on smart neighborhood
 23 functions; [52] is a good starting point for heuristics to map parallel programs on multi-
 24 core platforms. Self-adaptive Large Neighborhood Search has been successfully applied
 25 to scheduling problems [55] and to joint allocation and scheduling problems modeled
 26 via time interval variables [54] (see Section 3.1.2). Finally, some combinatorial prob-
 27 lems such as Capacitated Vehicle Routing with Time Windows are related to resource
 28 allocation and scheduling, but are left out of this survey: for a good starting point the
 29 interested reader may refer to [84].

30 The outline of the paper is as follows: Section 2 introduces the considered problem
 31 class and relevant variants; Section 3 discusses the main modeling techniques in OR and
 32 CP, including decomposition based models. Sections 4 and 5 are devoted to filtering
 33 algorithms and bounding rules; an overview of search strategies is given in Section 6,
 34 while solution methods for decomposition based approaches are discussed in Section 7.

38 2 Reference Problem Class

39 In this section, we provide a formal definition for the main problem class we take into
 40 account; relevant variants are discussed in Section 2.1.

41 We consider an allocation and scheduling problem defined over a set A of activities
 42 (say $\{a_0, a_1, \dots\}$), representing atomic operations. Activities have temporal dependen-
 43 cies, described by a directed acyclic graph $G = \langle A, E \rangle$, where they correspond to nodes
 44 and are connected by a set E of end-to-start precedence relations (a_i, a_j) . Without
 45 loss of generality, we assume there is a single source activity and a single sink activity.
 46 Following the OR literature, this will be referred to as *project graph*. We refer as s_i to
 47 the start time of activity a_i (i.e. the first time instant t where the activity is executing)
 48 and as e_i to the end time of a_i (i.e. the first time instant where the activity is not
 49

1 executing, after it has started), with $s_i, e_i \geq 0$; the notation \bar{s}/\bar{e} refers to the full array
 2 of start/end times.
 3

4 The problem defines a set R of resources r_k , with fixed capacity c_k over time (*re-*
 5 *newable resources*, introduced in [31]); those may represent, e.g., manpower, industrial
 6 machines or computing devices. We provide a formal description of allocation deci-
 7 sions by introducing binary variables x_{ik} , with $x_{ik} = 1$ if a_i is using r_k . Problem
 8 dependent constraints typically restrict the possible resource assignment choices, i.e.
 9 the allowed tuples for the full array of allocation variables \bar{x} . The exact amount of each
 10 requirement depends on the assignment decisions: formally, this is a function $rq_{ik}(\bar{x})$,
 11 with $rq_{ik}(\bar{x}) \geq 0$. Similarly, the activity durations are allocation dependent, i.e. each
 12 duration is a function $d_i(\bar{x})$, with $d_i(\bar{x}) \geq 0$. Note the approach allows duration and
 13 requirement values to depend on allocation choices involving multiple activities.

14 A solution of the problem is a full assignment of start times s_i and allocation
 15 variables x_{ik} such that neither temporal nor resource capacity restrictions are violated;
 16 finding a solution to the presented problem has polynomial complexity. A solution is
 17 optimal if it has (with no loss of generality) minimal value for a performance metric
 18 $F(\bar{x}, \bar{s}, \bar{e})$, depending in general both on allocation and scheduling choices; finding an
 19 optimal solution is an NP-hard problem (the proof follows from the one in [50]).

21 2.1 Relevant Problem Variants

22 Resource assignment and scheduling problems mainly arise in practical contexts; as a
 23 consequence, the basic structure outlined in the previous section is often complicated
 24 by a number of side constraints and unique features; rather than attempting a com-
 25 prehensive classification, we describe some of the main variants encountered in the
 26 literature and considered in this paper.

27 2.1.1 Resource Related Variants

28 Besides renewable resource, other resource types are found in the literature; here, we
 29 mention the main ones.

30 Non-renewable Resources

31 Non-renewable resources have a starting availability and are consumed throughout
 32 the whole scheduling horizon, until depleted: project budget is a good example. They
 33 were introduced in [31] and are typically considered by MRCPS approaches. If non-
 34 renewable resources are taken into account, finding a feasible resource assignment is
 35 an NP-complete problem [50] and the whole optimization process becomes consider-
 36 ably harder (see [76] for some comments). The so-called *doubly constrained resources*
 37 (coupling a fixed capacity over time and an overall usage limit) are considered in some
 38 works, but can be modeled by combining a renewable and a non-renewable resource.

39 Variable Capacity and Requirements over Time

40 Resources with time *varying capacity* are taken into account in [78, 76]. Work [8] shows
 41 a technique to turn variable capacities into constant ones, by introducing fictitious
 42 activities and fixing their position in the schedule by means of time windows or time
 43 lags (see Section 2.1.2); the method is described for the RCPSP, but can be easily
 44 generalized to take into account resource assignments. The important case of *resource*

1 *breaks* (e.g. vacation time, machine maintenance downtime) and break-interruptible
 2 activities (i.e. that can be preempted by breaks) is considered in [19] in an OR context
 3 for the MRCPSp and in [1,3] in CP. Time varying *resource requirements* are instead
 4 considered in [29] and in [69] in the context of alternative activities (see Section 3.1.2).
 5

6 *Time/Resource Trade-offs*

7 Some allocation and scheduling formulations take into account resource/time trade-
 8 offs, by allowing activities to assume a different duration, depending on the consumed
 9 amount of the selected resources. The MRCPSp is the most relevant example, since
 10 activity modes may specify different requirement values, rather than different mapping
 11 choices. In such a case, mapping decisions can no longer be formalized via the x_{ik} vari-
 12 ables. In this paper, we restrict to problems where time/resource trade-offs, provided
 13 they also feature resource mapping choices: this leaves out pure trade-off problems with
 14 fixed resource mapping, such as [27].
 15

16 *2.1.2 Temporal Constraint Related Variants*

17 Temporal constraints other than simple end-to-start precedence relations are very com-
 18 mon in the literature; here, we give a fairly complete overview of the reported cases.
 19

20 *Start/Start, Start/End, End/End Precedence Relations*

21 End-to-start relations can be generalized by introducing start-to-start, end-to-end and
 22 start-to-end precedence constraints. Those can be easily turned one into another, pro-
 23 vided activities have fixed durations, say d_i (for example $e_j \geq e_i$ becomes $s_j + d_j \geq e_i$).
 24 The transformation method is described in [30] for the RCPSP, but applies to the prob-
 25 lem from Section 2 provided $d_i(\bar{x})$ is constant, or whenever during search a full resource
 26 assignment becomes known. The multi-mode RCPSP with this kind of precedence re-
 27 lation is known as Generalized MRCPSp.
 28

29 As a relevant remark, if precedence relations other than end-to-start are taken
 30 into account, assigning the shortest possible duration to an activity may result in an
 31 increase of the schedule makespan [30]. This may prevent the application of several
 32 dominance rules (see Section 5.2.2).
 33

34 *Generalized Precedence Relations and Time Windows*

35 This case subsumes the previous one; additionally, minimal and maximal time lags (say
 36 δ^{min} and δ^{max}) label the precedence relations and constrain the time distance between
 37 the involved activities; formally in case of an end-to-start arc, the produced schedule
 38 must satisfy $\delta^{min} \leq e_j - s_i \leq \delta^{max}$. Minimal time lags enforce a minimum separation
 39 restriction, maximal time lags may model “best before” constraints, occurring e.g. in
 40 chemical and food industries. The Multi-mode RCPSP with this type of precedence
 41 constraint is known as MRCPSp with Generalized Precedence Relations (GPR). Time
 42 windows, constraining the execution of activities a_i between a release time rs_i and a
 43 deadline dl_i are equivalent to minimal/maximal time lags from the source node.
 44

45 For pure scheduling problems, a maximal time lag on an arc (a_i, a_j) can be con-
 46 verted into a negative minimal time lag on the complementary arc (a_j, a_i) ; the transfor-
 47 mation follows trivially by the inequality $e_j - s_i \leq \delta^{max}$ and is described in [8,18]. As
 48 a consequence, a graph with maximal time lags typically contains cycles; a temporally
 49 feasible assignment of start times can be still obtained via shortest path algorithms (as
 50

1 shown for the Simple Temporal Problem [26]), but finding a temporally and resource
 2 feasible schedule becomes NP-hard. Moreover, taking into account Generalized Prece-
 3 dence Relations may prevent the application of many dominance rules, as described in
 4 the previous paragraph.
 5

6 In allocation and scheduling problems, time lags may depend on the resource as-
 7 signment (e.g. mode dependent time lags in [73,38]); in the context of optimization for
 8 parallel programs on multi-core systems, allocation dependent minimal time lags are
 9 typically used to model communication costs depending on processor assignment decisions
 10 [51]. As a last relevant remark, observe that, as long as some assignment decisions
 11 have to be taken, allocation dependent durations are equivalent to start-to-end prece-
 12 dence constraints with allocation dependent time lags; as a consequence, most resource
 13 assignment and scheduling problem do behave as containing Generalized Precedence
 14 Relations, at least for part of the search process.

15 *Setup times*

16 Setup times are separation constraints between tasks mapped on the same resources,
 17 modeling, e.g., the time required to perform cleaning operations or to change tools
 18 before executing an activity a_i ; unlike allocation dependent minimal time lags, the
 19 involved activities need not to be connected by an arc in the project graph (i.e. their
 20 order is not a-priori fixed). If the setup time between activities a_i and a_j does not
 21 depend on the order they appear in the schedule, it may be incorporated in their
 22 execution time. Conversely, one speaks of *sequence dependent setup times*, a relevant
 23 case in practical applications; this is considered in [33] for scheduling with alternative
 24 resources, [10] for scheduling aircraft landings and [83] for a batching and scheduling
 25 problem from the chemical industry.
 26

27 *2.1.3 Objective Function Types*

28 *Time Based Objectives*

29 Time based objectives used in resource assignment and scheduling are extension of
 30 those used for pure scheduling problems. A time based cost function only indirectly
 31 depends on allocation choices, i.e. $F(\bar{x}, \bar{s}, \bar{e}) = F(\bar{s}, \bar{e})$. Makespan (overall completion
 32 time) minimization is the most frequently occurring objective in the literature; in
 33 this case we have $F(\bar{x}, \bar{s}, \bar{e}) = \max_i e_i$. Equivalently, assuming a_{n-1} is the single sink
 34 activity, one can state $F(\bar{x}, \bar{s}, \bar{e}) = e_{n-1}$.

35 In a real-world context, where oversubscribed problems with tight time-windows
 36 often arise, tardiness based costs are also very popular; for an activity a_i with deadline
 37 dl_i , the tardiness is the activity delay w.r.t the deadline; this is defined as $TD_i(e_i) =$
 38 $\max(0, e_i - dl_i)$. A tardiness cost function has the form $F(\bar{x}, \bar{s}, \bar{e}) = \sum_i w_i \cdot TD_i(e_i)$,
 39 where weights w_i are introduced for each activity to account for non uniform tardiness
 40 costs. This type of objective is considered in [41,54] from a methodological perspective
 41 and in [83,10,32,89] in an industrial context. Sometimes, there is some cost to be paid
 42 for early completions; earliness if formally defined as $ER_i(e_i) = \max(0, dl_i - e_i)$; ear-
 43 liness costs are considered in [83] (where they are due to inventory maintenance), in
 44 [10,32] (in case of early aircraft landing) and in [54] (from a methodological perspec-
 45 tive). A detailed list of possible problem objectives appears in [76]: most of them are
 46 time-based.

1 *Resource Based Objectives*

2 We refer as resource based objectives to cost functions directly depending only on
 3 allocation choices, i.e. $F(\bar{x}, \bar{s}, \bar{e}) = F(\bar{x})$. The main example are costs to be paid for
 4 assigning single activities to resources or for choosing a specific activity mode: those
 5 include processing costs [47], energy dissipation over multi-core systems [72], or as-
 6 signment costs in general [40, 82]. Costs associated to the assignment of activity *pairs*
 7 (similarly to those of a Quadratic Assignment Problem) appear instead in [16, 62] to
 8 model bus congestion in a multi-core system and are considerably more challenging
 9 from a computational standpoint.

10 *Time and Resource Based Objectives*

11 This class includes in first place any combination (e.g via weighted sum) of time-based
 12 and resource-based objectives. Moreover, some cost functions depend inherently on
 13 allocation and scheduling decisions; this is the case for setup costs, since they are
 14 defined for activities assigned to the same resource: those are considered in [33], in
 15 [83] (where they are due to machine cleaning operations) and in [72] (where they
 16 are associated to processor frequency switching). A second unusual example are the
 17 rescheduling costs taken into account in [25] in the context of reactive scheduling.

18 *Regular vs Non-Regular Objectives*

19 The notion of *regular performance measure* is introduced in [71] for the RCPSP and
 20 extended in [78] to the multiple mode case; a cost function is regular if its value may
 21 only improve by reducing the end time e_i of some activity (without changing the
 22 activity mode – i.e. the resource assignment). This is a *extremely* important case, since
 23 with regular objectives the set of optimal solutions always includes a *tight* schedule (see
 24 the detailed discussion on the Left-shift Rule in Section 5.2.2): this property is used
 25 to devise powerful pruning rules. Non-regular objectives are however very common in
 26 practice and include earliness costs, setup costs and all the mentioned resource-based
 27 objectives.

28 **3 Modeling Techniques**

29 Here, we present the main modeling techniques developed for the reference problem in
 30 CP (Section 3.1) and OR (Section 3.2); moreover, Section 3.3 discusses decomposition-
 31 based and hybrid approaches.

32 **3.1 Constraint Based Models**

33 *3.1.1 Using Classical Activities and Constraints*

34 It is possible to use classical CP techniques for pure scheduling problems to take into
 35 account resource assignment decisions. In Constraint-based Scheduling [3], activities
 36 are represented by introducing integer variables for every activity a_i . In detail we have
 37 S_i , representing the activity start time s_i and E_i , representing the activity end time
 38 e_i . In assignment and scheduling problems, it is customary to introduce an integer
 39 variable D_i representing the activity duration; then the constraint $E_i = S_i + D_i$ is
 40 enforced. Variables S_i , E_i and D_i are assumed to have convex domains; their bounds

```

1      min:  $F(\bar{\mathbf{X}}, \bar{\mathbf{S}}, \bar{\mathbf{E}})$ 
2      subject to:  $E_i = S_i + D_i \quad \forall a_i \in A$ 
3           $E_i \leq S_j \quad \forall (a_i, a_j) \in E$ 
4           $\text{cumulative}(\mathbf{S}, \mathbf{D}, RQ_{ik}, c_k) \quad \forall r_k \in R$ 
5           $D_i = d_i(\mathbf{X}) \quad \forall a_i \in A$ 
6           $RQ_{ik} = rq_{ik}(\mathbf{X})$ 
7          with:  $S_i, E_i, D_i \in \{0, \dots, eoh\}$ 
8           $X_{ik} \in \{0, 1\}$ 
9

```

Fig. 1 A CP model for non-preemptive resource allocation and scheduling, with no side constraints

are referred to by means of conventional names: $\min(S_i)$ is the Earliest Start Time – $EST(a_i)$ – and $\max(S_i)$ is the Latest Start Time – $LST(a_i)$; the Earliest End Time and Latest End Time – $EET(a_i)$ and $LET(a_i)$ – are defined analogously for the E_i variable.

End-to-start precedence relations are modeled as linear constraints $E_j \leq S_j$. Generalized Precedence Relations can be taken into account by posting constraints $\delta^{min} \leq E_i - S_i \leq \delta^{max}$, according to the definition from Section 2.1.2. Enforcing consistency for the resulting network is known as Simple Temporal Problem [26] and can be done via propagation of the inequality constraints; however, detecting an infeasibility with this method takes in general time proportional to the largest S_i/E_i domain. Alternatively, precedence constraint can be explicitly stated (e.g. by means of global constraints) and one may use a shortest path algorithm for graphs with negative weights, such as Bellmann-Ford (with complexity $O(|A||E|)$, where $|A|$ is the number of activities and $|E|$ is the number of arcs).

Allocation decisions can be modeled [2] by introducing a binary variable X_{ik} for each x_{ik} ; duration variables are linked to allocation decisions by a constraint $D_i = d_i(\mathbf{X})$, where \mathbf{X} represents the whole set of X_i variables and the function $d_i(\cdot)$ is the one from Section 2. Side constraints usually restrict the possible resource assignments. Resource restrictions are enforced via the *cumulative* constraint (see [3] for a reference); the effect of resource assignments can be modeled by using variable resource requirements: in detail, the amount by which activity a_i requires resource r_k is modeled via an auxiliary variable RQ_{ik} , depending on allocation decisions, i.e. $RQ_{ik} = rq_{ik}(\mathbf{X})$, where $rq_{ik}(\cdot)$ is the function from Section 2. Classical filtering for the cumulative constraint can be used in this case by assuming RQ_{ik} has minimal value.

Figure 1 shows a basic CP model for resource allocation and scheduling with no side constraint; the cost function F is as from Section 2 and eoh denotes the largest considered time instant (End Of Horizon). Non-renewable resources are straightforward to model with this approach and result into a set knapsack constraints on X_{ik} variables. Similarly, time-based objective and assignment costs are easy to deal with. Sequence dependent setup times are usually handled by introducing sequencing variables B_{ij} such that $B_{ij} = 1$ iff a_i precedes a_j , i.e. we have $B_{ij} = 1 \Leftrightarrow E_j \geq S_i + \delta_{ij}$, where δ_{ij} is the setup time. From a modeling perspective, there is hardly any allocation and scheduling problem beyond the reach of this approach; unfortunately, resource and temporal propagation tend to be very weak when requirement or duration variables are unbound, so that more structured approaches are in practice needed.

1 3.1.2 Specific Constraint Based Modeling Techniques
 2
 3
 4

5 *Alternative Resources*
 6
 7
 8
 9
 10
 11
 12
 13
 14
 15

This approach models assignment choices by allowing an activity a_i to require exactly one out of a set, say $R' \subseteq R$, of alternative resources [33,65]; formally, relation $\sum_{r_k \in R'} x_{ik} = 1$ must hold. Typically, the whole resource set R is partitioned into subsets R^0, R^1, \dots of alternative resources. A fixed requirement (say ρ) is specified for the whole set, so that the corresponding $r_{qik}(x)$ function is $\rho \cdot x_{ik}$. We say an alternative resource set R' is independent if a decision on R' has no impact on any other assignment choice, i.e. if there is no constraint connecting x_{ik} variables with $r_k \in R'$ to other variables in \bar{x} . From a modeling perspective, alternative resources represent a restriction compared to the approach in the previous section; however, they allow more powerful filtering, taking advantage of the fixed requirement value and the XOR relation between x_{ik} variables.

As a relevant special case, an independent set R' of n identical unary resources is equivalent to a single resource with capacity n (multi-capacity resource, see [64]); this representation results into huge search time savings whenever applicable. Unfortunately, side constraints or assignment dependent durations may make the resources non-identical and invalidate the equivalence; setup times have the same effect, as they require to know the specific resource assignment for each activity.

22
 23 *Disjunctive Temporal Problem with Finite Domain Constraints*
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35
 36
 37
 38

Assignment and scheduling problems over unary capacity resources can be modeled within the framework of the Disjunctive Temporal Problem with Finite Domain Constraints (DTP_{FD}) by Moffitt, Peintner and Pollack [63]; a DTP_{FD} defines a network of time points, each with an associated temporal variable T_i and a finite domain variable Y_i ; time points can represent activity start/end events. It is possible to post between pairs of time points a *disjunction* of linear inequalities in the form $T_i - T_j \leq DB(Y_i, Y_j)$, with the value of the bounding function $DB(Y_i, Y_j)$ depending on the finite domain variables. This approach exposes in depth the temporal structure of the problem, providing support for stronger filtering between time and assignment variables (see Section 4.1.1). Variables Y_i can be linked to assignment variables x_{ik} so that two activities are required not to overlap if they are processed by the same unary resource. The approach naturally handles setup times, by connecting sequencing variables (e.g. the B_{ij} we used before) to Y_i variables by means of chaining constraints. Moreover, the DTP_{FD} can be used to model time/resource trade-offs (see Section 2.1.1).

Non-unary resources can in principle be handled by detecting possible resource conflicts at search time and consequently posting new disjunctions, similarly to what is done in many Precedence Constraint Posting approaches (see Section 3.3); this is however not reported in the literature.

43
 44 *Alternative Activities*
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

The so-called Alternative Activities have been introduced by Beck and Fox in [11,12] and can be used to model resource assignment decisions and their impact on durations and requirements. The project graph is extended by including so-called *XOR nodes*, marking the start (and the end) of alternative blocks. Formally, each activity is assigned an execution variable ex_i , such that $ex_i = 0$ if activity a_i does not execute, $ex_i = 1$ in case a_i executes and $ex_i \in \{0, 1\}$ as long as it is undecided. XOR nodes mark the

beginning and the end of alternative subgraphs. Formally, successors and predecessors of a XOR node are mutually exclusive (i.e. $\sum_i ex_i = 1$); blocks of alternative subgraphs are required to be nested. A block with a single activity on each branch can be used to represent different recipes to execute the same logical activity (i.e. combinations of duration and resource requirements), equivalently to modes in the MRCPSP. Alternatives may differ by the required resources or by the required amount, so that this approach can be used to model time/resource trade-offs (unlike alternative resources). Alternative paths between nested XOR nodes may contain generic subgraphs rather than single activities: this allows one to model alternative process plans, making this approach more expressive than the MRCPSP.

Alternative activities and plans are considered (from a temporal perspective only) in Temporal Networks with Alternatives (TNA), introduced by Barták, Cepek and Surynek in [5]. A TNA is very similar to a graph with XOR nodes, the approach is in principle not restricted to nested alternative blocks, allowing one to express more complex assignment constraints; unfortunately, work [5] proves that completing a partial assignment of execution variables (e.g. finding a feasible resource assignment) is NP-complete, while the problem is polynomial in nested networks² [4].

In general, exposing assignment choices as alternative activities provides more information to filtering algorithms (see Section 4.1.1). As a main drawback, the project graph and the model get bigger; if Generalized Precedence Relations are not considered, modeling assignment decision over independent sets of alternative resources requires an exponential number of activities (i.e. one for each combination of resource assignments); conversely, by using GPRs it is possible to avoid the combinatorial blow-up by building an alternative block for each independent set of resources and synchronizing their opening/closing XOR nodes using precedence constraints with time lags.

Optional activities and Interval Variables

The so-called optional activities are similar to alternative activities since they have an associated execution variable ex_i ; however, they are not necessarily part of an alternative block. They are introduced by Le Pape in [58], where the activity representation of ILOG-Scheduler is extended, so that a 0 durations denotes a non-executing activity. Optional activities requiring different resources can be used to model complex resource assignments decisions, by properly constraining their execution variables; the approach is more expressive than alternative activities, but may lead to weaker propagation (see Section 4.3).

Building over ideas developed for alternative and optional activities, in [56], the constraint engine has been extended by Laborie and Rogerie to handle optional activities as first class variables: those are referred to as *time-interval variables*. A time interval variable \underline{A}_i has values in the domain $\perp \cup \{[s, e] \mid s, e \in \mathbb{Z}, s \leq e\}$; namely, either the variable is *non-executed* ($\underline{A}_i = \perp$) or its value is a half-open interval $[s, e]$ with integer bounds. Non-executed variables have no effect on the constraints they are involved in. Time interval variables can be connected by generalized precedence constraints, or can be organized in explicit alternative and hierarchical blocks; each block corresponds to a macro-interval \underline{A}_i , spanning over a set of (possibly alternative) sub-intervals \underline{A}_j .

So-called execution constraints $exec(\underline{A}_i)$ force a variable to be executed and can be aggregated into unary or binary logical expressions (e.g $exec(\underline{A}_i) \Rightarrow exec(\underline{A}_j)$). Resource

² We conjecture the result may hold for slightly more general graph structures, such as those considered in [62].

$$\begin{aligned}
& \min F(\bar{\mathbf{E}}) && (1) \\
& \text{s.t. } \sum_{\mu_h \in M_i} \sum_{t=0}^{eoh} \mathbf{E}_{i,h,t} = 1 \quad \forall a_i \in A && (2) \\
& \sum_{\mu_h \in M_i} \sum_{t=0}^{eoh} t \cdot \mathbf{E}_{i,h,t} \leq \sum_{\mu_h \in M_j} \sum_{t=0}^{eoh} (t - d_j(\mu_h)) \cdot \mathbf{E}_{j,h,t} \quad \forall (a_i, a_j) \in E && (3) \\
& \sum_{a_i \in A} \left[\sum_{\mu_h \in M_i} rq_{i,k}(\mu_h) \sum_{t=t_0}^{t_0+d_i(\mu_h)-1} \mathbf{E}_{i,h,t} \right] \leq c_k \quad \forall r_k \in R, \forall t_0 = 0 \dots eoh && (4) \\
& \mathbf{E}_{i,h,t} \in \{0, 1\} \quad \forall a_i \in A, \forall \mu_h \in M_i, \forall t = 0 \dots eoh && (5)
\end{aligned}$$

Fig. 2 A Time Indexed model for non-preemptive MRCPSP

constraints are taken into account in [57], where time-intervals are associated to step functions (referred to as cumul functions), representing the resource usage profiles. The approach is showcased in [54] on three practical examples.

Complex resource allocation and scheduling problems can be modeled by encoding resource assignment as execution decisions for interval variables; the idea of exposing different execution recipes as different activities, explicit specification of alternative blocks and the availability of GPRs provide strong support for filtering algorithms; the approach is currently adopted by IBM-ILOG CP Optimizer and Google or-tools.

3.2 Mixed Integer Linear Programming models

Most of the OR literature about resource allocation and scheduling is related to the Multi-mode Resource Constrained Scheduling Problem (MRCPSP), first introduced in [31]. In the MRCPSP, each activity a_i can be executed in one out of a set of possible modes M_i ; each mode $\mu_h \in M_i$ represents an alternative way to execute the activity and specifies a set of resource requirements and a duration value. The MRCPSP does not strictly fit the notation introduced in Section 2, but can be taken into account by introducing a mode variable m_i for each activity a_i and making the requirements and duration functions mode dependent, i.e. $rq_{i,k}(\bar{x}) \rightarrow rq_{i,k}(m_i)$, $d_i(\bar{x}) \rightarrow d_i(m_i)$.

The multi-mode formalism allows one to model time/resource trade-offs, since different modes may require the same resources in different amounts; similarly to alternative activities and time interval variables, modeling assignment decisions over independent set of alternative resources requires an exponential number of modes, unless GPRs are supported. The main Mixed Integer Linear Programming (MILP) models for the MRCPSP can be classified into *time indexed* and *disjunctive* and are described in the following.

3.2.1 Time Indexed Model

In time indexed models binary variables $\mathbf{E}_{i,h,t}$ are introduced to denote whether an activity a_i is scheduled to finish at time t in mode μ_h ; in Figure 2, we report the model introduced by Talbot in [81]. Constraints (2) require each activity to be finished by the end of horizon. Constraints (3) enforce end-to-start precedence relations and

1 constraints (4) resource capacity restrictions. The cost function has been re-formulated
 2 to take into account the different decisions variables.
 3

4 Time indexed models allow resource capacity restrictions to be formulated as linear
 5 constraints; as a drawback, due to the use of a discrete time representation, the number
 6 of variables depends on the length of the horizon, resulting in limited scalability. In
 7 an attempt to address the issue, work [90] proposes two alternative formulations, but
 8 reports no substantial improvement. With time indexed models it is straightforward
 9 to represent non-convex temporal domains, thus allowing to easily handle scenarios
 10 requiring disjoint time windows, e.g. resource breaks with non-interruptible activities.
 11 Variable resource capacities over time are similarly easy to deal with (see [91]). The
 12 linear relaxation provided by a time indexed model is usually stronger than that of
 13 disjunctive one, in particular in case of tight resource constraints.

14 *Disjunctive Model*

15 In a disjunctive model, a start variable s_i is introduced for each activity a_i ; mode
 16 assignments are represented by variables M_{ih} , such that $M_{ih} = 1$ iff activity a_i is executed
 17 in mode μ_h . A complete model (a slight elaboration over the one by Heilmann in [38])
 18 is shown in Figure 3. Constraints (7) represent end-to-start precedence relations and
 19 Constraints (9) force a mode to be assigned to each activity. The notation $A(\bar{s}, t)$ refers
 20 to the set of tasks executing at time t , i.e. such that $s_i \leq t < s_i + \sum_{\mu_h \in M_i} d_i(\mu_h) \cdot M_{ih}$.
 21 The cost function has been re-formulated to take into account the different decisions
 22 variables.
 23

24 Disjunctive models require a smaller number of decision variables compared to time
 25 indexed ones; however, Constraints (8) do not have a simple linear representation. This
 26 can be provided by preventing the overlapping execution of conflicting activities and
 27 requires: (1) to identify the sets of activities possibly causing an overusage (e.g. the
 28 Minimal Forbidden Sets described in Sections 6.1 and 6.2); (2) to add constraints
 29 forcing the disjoint execution of some of the activities in the sets. Such a method has
 30 however two strong disadvantages: (1) the number of Minimal Forbidden Sets (and
 31 therefore of additional constraints) is in general exponential in the size of the graph³;
 32 (2) the disjunction constraints make use of a big-M linearization and lead to a poor
 33 relaxation based bound (considerably worse than time-indexed models). As a matter
 34 of fact, all approaches based on disjunctive models such as [38] rely on specific search
 35 strategies to take care of resource constraints (see Section 6).

36 As an alternative, [73] proposes a more compact model where resource allocation
 37 and scheduling is formulated as a constrained rectangular packing problem. Unfor-
 38 tunately, since rectangular packing and renewable resource restrictions are not fully
 39 equivalent [13], the method may miss feasible/optimal solutions.
 40

41 3.3 Hybrid and Decomposition Based Approaches

42 Resource assignment and scheduling problems result from the composition of a pure
 43 assignment and a pure scheduling component; this hybrid nature can be exploited to
 44 split the overall problem into distinct stages with separated (but dependent) variables.
 45 Typically, one is left with a resource assignment problem (so-called master problem)

46 3 It is quadratic if only unary resources are taken into account, making this approach fairly
 47 popular in such context.
 48

1	$\min F(\bar{\mathbf{S}}, \bar{\mathbf{M}})$	(6)
2	s.t. $\mathbf{S}_j - \mathbf{S}_i \geq \sum_{\mu_h \in M_i} d_i(\mu_h) \cdot \mathbf{M}_{ih}$	(7)
3	$\forall (a_i, a_j) \in E$	
4	$\sum_{a_i \in A(\bar{\mathbf{S}}, t)} \sum_{\mu_h \in M_i} r q_{i,k}(\mu_h) \cdot \mathbf{M}_{ih} \leq c_k$	(8)
5	$\forall r_k \in R, \forall t = 0 \dots eoh$	
6	$\sum_{\mu_h \in M_i} \mathbf{M}_{ih} = 1 \quad \forall a_i \in A$	(9)
7	$\mathbf{S}_i \geq 0 \quad \forall a_i \in A$	(10)
8	$\mathbf{M}_{ih} \in \{0, 1\} \quad \forall a_i \in A, \forall \mu_h \in M_i$	(11)

Fig. 3 A Disjunctive model for non-preemptive MRCPSP

and a pure scheduling problem (so-called subproblem), given the fixed allocation provided by the first stage.

Such a decomposition has a number of advantages: (1) models for both the stages are typically much smaller and have a cleaner structure, making their solution more efficient: for instance, it is possible to remove big-M expressions from a MILP model, improving the linear relaxation [47]. (2) Heterogeneous techniques can be used to solve the stage they are best for, e.g. MILP for the assignment and CP for scheduling, obtaining so-called hybrid approaches. (3) “Pure” problems can be tackled with a much better established pool of techniques, such as those in [20, 3]. (4) In case of resource-based objectives (see Section 2.1.3), the cost function depends only on the assignment stage, with the scheduling one acting as a feasibility check, e.g. [47, 40, 16, 62]. (5) Finally a clever decomposition may allow the scheduling stage to be split into a collection of independent subproblems, making its solution even simpler; this is relatively easy if the Project Graph contains no precedence relations or in case only setup times are specified, but may be impossible if no special precedence constraint structure is assumed. This approach is adopted e.g. in [40, 21].

Alternatively, a problem decomposition can be obtained by separating constraints rather than problem variables. In such a case the master is an optimization problem, possibly formulated with some technique with no support for resource constraints; the subproblem is either a consistency check or a very simple feasibility problem. The approach is discussed from a general perspective in [39], and in [82] in the context of the Branch-and-Check framework. Typically, resource capacity constraints can be considered only in the subproblem, leaving temporal and assignment constraints in the master; this is the key idea behind many search schemes adopted in Precedence Constraint Posting, based on the identification of resource conflicts in the current schedule (see Sections 6.1 and 6.2) and their resolution by adding cuts; the approach has been devised for pure scheduling problems (e.g. [70, 53]) and applied to allocation and scheduling problems in [38, 12].

The obvious drawback of decomposition is the loss of valuable information due to the decoupling; in order to improve the solution quality, the two stages are solved iteratively in an interacting fashion: in particular, (1) cuts are injected in the assignment problem whenever the scheduling problem is solved (either successfully or not); moreover (2) in all the mentioned references the assignment problem contains some relaxation of the scheduling problem constraints (so-called subproblem relaxation). Properly engineering the interaction between master and sub-problem is essential for

the effectiveness of a decomposition based approach: systematic, optimal, methods to achieve this goal are formalized in the Branch-and-Check [82] framework and (most relevantly) in Logic-based Benders' Decomposition (LBD) formally defined by Hooker and Ottosson in [42]; both the approaches are discussed in Section 7.

A recent paper [83] tackles an allocation, batching and scheduling problem in the chemical industry introducing *temporal* decomposition: the time horizon is split into time buckets; an LBD like approach is used to solve an allocation and scheduling problem for each bucket, with resource assignment acting as the master problem; the schedule for time bucket i is built by extending that of bucket $i - 1$; solved buckets are never revised. Finally, it is possible to devise hybrid approaches not relying on decomposition: work [47] reports a double CP/MILP model for an assignment and scheduling problem; where branching is essentially based on the CP representation, while the MILP model is used to obtain bounds via its linear relaxation.

4 Propagation

In the context of resource assignment and scheduling problems, the key difficulty addressed by propagation techniques is to provide the tightest possible interaction between allocation and scheduling choices; in the following, we overview existing propagation techniques, roughly categorized in *temporal*, *logical* and *resource* filtering, depending on the aspect they mainly focus on. It is interesting to see how most of the filtering techniques for assignment and scheduling problems are in fact relatively simple extension of the corresponding methods for pure scheduling: we found this stimulating rather than disappointing, as it leaves room for possible future improvements. Finally note that, with decomposition methods, classical scheduling constraints can be used in each subproblem, while the problem of the interaction between resource assignment and scheduling translates to that of the interaction between the resulting components (addressed in Section 7).

4.1 Temporal Reasoning

If only end-to-start precedence relations are defined, temporal constraint propagation can be done via the classical Critical Path Method [48] and lower bounds on the time windows (i.e. LSTs) can be obtained by assuming the durations are fixed to the lowest value and computing shortest paths from the source node. Temporal reasoning and time window determination have an important role in reducing the size of time-indexed models (see Section 3.2.1) and the quality of their linear relaxation.

In case maximal time lags or precedence constraints other than end-to-start are considered, minimum durations do not necessarily lead to valid LSTs; this happens if a backward arc is part of the critical path as discussed in [30]. Correct time windows can be computed by relying on the equivalence between durations and min/max time lags (see Section 2.1.2), i.e. $\min(d_i(\bar{x})) \leq e_i - s_i \leq \max(d_i(\bar{x}))$: enforcing bound consistency on the resulting inequality constraints or performing shortest/longest path computation on the equivalent Simple Temporal Problem leads to valid EST/LET values. Obviously, no issue exists at no duration is allocation dependent; in general, however, when complex precedence constraints are taken into account, one should be

1 aware that shortest durations do not necessarily result in the shortest possible schedule;
 2 this prevents the application of some dominance rules (see Section 5.2.2).
 3
 4

5 4.1.1 Temporal Reasoning with Alternative Resources and Activities

6 Alternative and optional activities can be tackled by including the execution variables
 7 $ex_i \in \{0, 1\}$ in the precedence constraint expression, so that an end-to-start constraint
 8 becomes e.g. $ex_i \times ex_j \times s_j \geq ex_i \times ex_j \times e_i$, as in [6]. The resulting propagation is
 9 however usually weak: this is quite expected since no specific constraint structure is
 10 assumed for the ex_i variables.
 11

12 Stronger propagation can be obtained for nested alternative blocks [12, 6], for al-
 13 ternative groups of time-interval variables [56] and for alternative resources: the basic
 14 technique used to perform filtering is the same in all cases and is described here in
 15 the context of alternative resources. In detail, let a_i be an activity requiring ρ units
 16 of one out of a set R' of m alternative resources, propagation can be performed as if
 17 a_i was split in m alternative activities a_i^0, \dots, a_i^{m-1} , each requiring ρ units of a specific
 18 resource in $r_k \in R'$. Then temporal reasoning maintains the constructive disjunction
 19 between the time window of the alternative sub-activities (see [33, 64]):
 20

$$\begin{aligned} EST(a_i) &= \min_{r_k \in R'} \{EST(a_i^k)\} & LST(a_i) &= \max_{r_k \in R'} \{LST(a_i^k)\} \\ EET(a_i) &= \min_{r_k \in R'} \{EET(a_i^k)\} & LET(a_i) &= \max_{r_k \in R'} \{LET(a_i^k)\} \end{aligned}$$

21 classical filtering methods for temporal constraint can be used to prune the time win-
 22 dows of the sub-activities; whenever the domain of a start time variable gets empty,
 23 the corresponding resource r_k is removed from the set of possible assignment choices,
 24 or the corresponding ex_i^k variable is set to 0 in case of alternative activities or time
 25 interval variables.

26 A similar approach is adopted by the least commitment method for the DTP_{FD}
 27 described in [63]; here, temporal constraints with disjunctive bounds in the DTP_{FD}
 28 are associated with min/max time lags by convexification (i.e. by assuming the loos-
 29 est bounds $DB(Y_i, Y_j)$ allowed by the finite domain variables); then relaxed, convex,
 30 time widow are computed for each node. Whenever the time windows are tightened,
 31 incoherent Y_i values are ruled out.

32 As a last, very interesting case, the time-interval variable approach described in
 33 [56], allows *binary* logical constraints to be specified on the ex_i variables; those are
 34 encoded as an implication graph, allowing some joint logical and temporal propagation
 35 to be performed; in detail, assuming an arc (a_i, a_j) exists, whenever the logical network
 36 can infer the relation $ex_i \Rightarrow ex_j$ the arc can propagate the conditional bounds from s_j
 37 to e_i . Similarly, if the relation $ex_j \Rightarrow ex_i$ can be inferred, then the other half of the
 38 propagation can be performed. This allows propagation of temporal bounds even if the
 39 execution status is not yet fixed.

40 4.2 Logical Filtering

41 This section mainly discusses the approach described in [7] for Temporal Network
 42 with Alternatives, that improves propagation by adding redundant side constraints
 43 on the execution variables. Such redundant constraints are based on the identification
 44

1
2
3
4
5
6
7
8
9
of classes of “equivalent” activities, inferred from the network structure. Equivalent
activities share the same execution variable: explicitly stating this information sub-
stantially improves propagation over resource assignment variables. Despite detecting
all equivalent nodes in a network is NP-hard, the paper proposes three polynomial
complexity methods to identify part of them; the most effective technique relies on
the detection of nested substructures, in which case the entry and exit node can be
considered equivalent. A second relevant example is the implication graph in [56], that
may allow stronger propagation compared to local consistency.

10
11
12
4.3 Resource Reasoning

13
14
The earliest edge-finding algorithm for resource assignment and scheduling problems
15
is discussed in [2]: optional activities are tackled by assuming non-executed activities
16
to have zero duration. Alternative resource assignments are dealt with by conditioning
17
the resource requirement with the demand variables x_{ik} from Section 3.1; the resulting
18
propagation is however quite weak.

19
Optional activities modeled via execution variables are first taken into account
20
into edge-finding reasoning in [12]; the algorithm is based on two simple ideas: (1) an
21
optional activity a_i should not be taken into account when computing bounds on any
22
other activity a_j ; (2) if the time window of a_i gets empty due to resource propagation,
23
the value 1 should be removed from the execution variable ex_i (i.e. the activity is
24
deemed non-executable). This kind of reasoning is the same encountered for temporal
25
propagation in Section 4.1.1 and is the key idea behind many other filtering algorithms,
26
extending classical propagators for resource constraints to optional activities. This is
27
done in [88] in the context of efficient edge finding with Θ -trees; in [33] for alternative
28
resources, where the execution variables are replaced by resource assignment choices; in
29
[87] for mixed timetable/edge-finding; in [86] within an interesting method to filter the
30
maximal resource usage and maximal duration; in [69], where activities with negative
31
resource consumption (i.e. resource producers) are taken into account.

32
The opposite process (deducing the necessary execution of an activity) is inherently
33
more difficult: typically, the most effective way to get to such a conclusion is propagating
34
side constraints on the resource allocation (or the execution) variables. As an exception,
35
in [69] it is shown how to deduce the necessary execution of activities with *negative*
36
resource consumption.

37
For a set R' of alternative resources, complementary propagation is possible by
38
introducing a redundant resource [64] with capacity $\sum_{r_k \in R'} c_k$ and applying classical
39
filtering algorithms; combined with the previously discussed methods, the approach
40
provides improved propagation as long as the resource assignment is undecided.

41
42
5 Lower Bounds and Bounding Rules

43
Research efforts in the context of the Multi-mode RCPSP have focused on bounding
44
and dominance rules. Bounding rules refer to inferring necessary restrictions on time
45
windows, mode assignments and scheduling decisions; they are closely related to filter-
46
ing in CP. Dominance rules are based on some proof that the set of optimal solutions
47
must include a schedule with specific properties, which are then used to narrow the
48
search space. As an exception, some dominance rules exploit knowledge on past search
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 to deduce that no better schedule can be reached from the current search node: they
 2 can be considered as a form of no-good learning. The most relevant bounding and dom-
 3 inance rules are presented in Section 5.2. Conversely, lower bounds for the MRCPSp
 4 tend to be not so effective: they are discussed in Section 5.1.

5

6 5.1 Lower Bounds

7

8 Lower bounds on the optimal problem cost are a fundamental component of many
 9 search strategies (such as branch & bound, branch & cut, branch & price); in CP, a
 10 lower bound can be encapsulated in a global constraint, improving the effectiveness
 11 on optimization problems and providing access to useful information, such as reduced
 12 costs [34].

13 Effective lower bounds for resource allocation and scheduling problem are difficult
 14 to obtain. Time indexed and disjunctive models for the MRCPSp provide ready to use
 15 bounds via their linear relaxation; unfortunately, those are not so tight. The relaxation
 16 of time-indexed models provides the best results, due to lack of linearized constraints
 17 via big-Ms. To the authors knowledge, [91] is the only work directly exploiting a bound
 18 obtained by the relaxation of a time-indexed model; the bounding technique is employed
 19 in two stages: in a preprocessing step to tighten time windows and reduce the model
 20 size, then during the solution process to fathom search nodes. The tightness of the
 21 linear relaxation can be improved by adding redundant cuts: this is done in [91] for the
 22 MRCPSp and in [10] in the context of scheduling aircraft landings.

23 In case of regular performance measures (see Section 2.1.3) a straightforward lower
 24 bound on the optimal cost can be obtained by ignoring resource constraints and start-
 25 ing each non-scheduled activity at the Earliest Start Time, obtained via longest path
 26 computation. Interestingly, this is the most widely employed bound for resource assign-
 27 ment and scheduling problems, since it is both easy to obtain and reasonably effective.
 28 In case of non-regular objectives, things get more complicated since an earliest start
 29 schedule does not necessarily correspond to a better cost: for earliness and tardiness
 30 objectives, [32] develops a modified network simplex algorithm which can be used to
 31 obtain in polynomial time an optimal schedule (with no resource constraint). Lower
 32 bounds for sequence dependent setup times are considered in [33], while work [25]
 33 discusses lower bounds for rescheduling costs.

34 An interesting technique exploited in [91] and [25] to strengthen a lower bound
 35 consists in performing truncated tree search with a maximum depth; the weakest bound
 36 on the tree frontier is valid for the root node.

37

38 5.2 Bounding Rules

39

40 Bounding rules are tree reduction techniques that check if the current search node can
 41 be preventively fathomed. Unlike filtering algorithms, bounding rules are executed as
 42 part of the search method and are not attached to a constraint; as a consequence, the
 43 rule formulation is tailored on a specific branching scheme (see Section 6.2), despite the
 44 main underlying idea usually has broader applicability. There is no general coordination
 45 mechanism (such as propagation), so that the combination of different rules is up to
 46 the algorithm designer. From a CP perspective, bounding rules may serve as a basis
 47 for the development of filtering algorithms.

1 5.2.1 Feasibility Based Rules
 2
 3
 4

5 *Non-renewable Resource Rule*

6 This rule appears in [37], in [77] (with the name “non-renewable resource consump-
 7 tion”) and in [25] (as “resource infeasibility rule”). The rule considers each *non-*
 8 *renewable* resource r_k : if scheduling each currently unscheduled activity in the mode
 9 with the lowest request for r_k would exceed its capacity c_k , then the current partial
 10 schedule cannot be feasibly completed. The rule is given a very efficient static formula-
 11 tion in [76] (where it is referred to as “input data adjustment”). This kind of reasoning
 12 is subsumed in CP by usual constraint propagation on the resource capacity constraints
 13 and the assignment variables.

14 *Non-delayability and Immediate Selection*

15 Those two rules, respectively reported in [77, 76, 37, 18] and [38] are based on the prin-
 16 ciple that a forced scheduling choice should be immediately performed. The general rule
 17 states that, if the current set of scheduling options contains an activity with no other
 18 chance to be scheduled, the choice should be immediately performed. The mentioned
 19 works contain specific adaptations to the different scheduling schemes.

20 *Single Enumeration Rule*

21 The single enumeration rule, introduced in [78] and further applied and refined in [76,
 22 37] is a type of dynamic symmetry breaking constraint for precedence tree branching
 23 (see Section 6.2). The rule targets two activities a' , a'' , scheduled in two subsequent
 24 search steps i and $i+1$ in mode m' and m'' ; if their assigned start times do not depend
 25 on which activity is scheduled at step i and which one at $i+1$, then only one sequence
 26 needs to be considered.

27
 28 5.2.2 Dominance Rules
 29
 30

31 Dominance rules are symmetry breaking constraints for optimization problems; they
 32 are based on the observation that the set of optimal solutions necessarily contains a
 33 schedule with specific properties; one can therefore focus on the generation of that
 34 specific schedule, considerably narrowing the search space.

35 Dominance rules come with specific applicability assumptions: in particular, the
 36 cost function is typically assumed to be regular (see Section 2.1.3); the applicability of
 37 most rules to non-regular objectives is scarcely discussed in the literature. All presented
 38 dominance rules target a search process where a partial schedule is built by assigning
 39 a start time to unscheduled activities, proceeding in chronological order.

40
 41 *Left-shift Rule*

42 This extremely important class of dominance rules is based on a property of regular
 43 objective functions (see Section 2.1.3), and on the concept of left-shift (discussed in
 44 details in [79]). A left-shift is an operation on a single activity a_i within a feasible
 45 schedule S , deriving a feasible schedule S' , such that $e_i(S') < e_i(S)$ (where $e_i(S)$ is
 46 the end time of a_i in S) and no other schedule modification occurs; a left shift of
 47 exactly one time unit is called *local left shift*. A multi-mode left shift [77] is a left shift
 48 of a_i where the activity is allowed to change mode (e.g. to be assigned to a different
 49 resource). Then, a schedule is *active* if it is feasible and no activity can be left-shifted;

50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

a schedule is *tight* if it is feasible and no multi-mode left-shift can be performed. In case of regular performance measures (see Section 2.1.3), the set of optimal schedules is guaranteed to contain a tight schedule (an active schedule in the case of the RCPSP).

A pruning rule can be devised based on those properties; the general version of this left-shift rule states that a partial schedule in which an activity a_i can be left-shifted without violating the precedence and the resource constraints needs not to be completed (as it is dominated by another active or tight schedule). Several variants of the rule exist: they are employed in [76, 25, 37, 77]; brief and effective descriptions are reported in [18].

12 *Multi-mode Rule*

This rule (employed in [76, 37, 18, 77, 24]) is based on the so-called mode-reduction operation. A mode reduction [77] of an activity a_i within a feasible schedule is an operation changing the mode of a_i to one with shorter duration, without changing its finish time and without violating the constraints or changing the modes or finish times of the other activities. A schedule is called *mode-minimal* if no mode reduction can be performed. If there is an optimal schedule for a given instance, then there is an optimal schedule which is *both tight and mode-minimal*; some care must be observed with maximal time lags and precedence constraints other than end-to-start (see Sections 2.1.2 and 4.1). Note there may be tight schedules which are not mode-minimal, and mode-minimal schedules which are not tight (for an example see [77]). The rule states that, if a mode reduction can be performed on an activity a_i with e_i equal to the current scheduling time, then the current partial schedule needs not to be completed.

26 *Order Swap Rule*

An order swap [37, 18] is an operation on a feasible schedule targeting two activities a_i, a_j with $j > i$, such that a_i, a_j are scheduled in sequence, i.e. $e_i = s_j$. The order swap consists in an exchange of the start time of the two activities, with no violation of precedence or resource constraint; changing the mode of any activity or the start time of any activity other than a_i and a_j is not allowed.

A schedule where no order swap can be performed is called *order monotonous*. If the order swap does not affect the objective function value (this is the case e.g. for the makespan, but not for tardiness costs), the set of order monotonous schedules is guaranteed to contain an optimal schedule. Therefore, before an activity a_i is scheduled at time t , if an order swap is allowed with any scheduled activity having end time t , then the current search node needs not to be further extended. Note the order swap rule should not be confounded with the single enumeration one, since the latter does not require the activities to form a sequence in the schedule.

42 *5.2.3 Static Bounding/Dominance Rules*

44 Static bounding rules are introduced in [77] and used in most of the exact approaches
45 for the MRCPS developed later on; they are applied prior to the beginning of search
46 and consist in the removal of *non-executable modes*, *inefficient modes* and *redundant*
47 *resources*. The rule application is iterative, until a fix-point is reached, making them
48 very similar to constraint propagation.

49 In detail, a mode μ_h for activity a_i is defined as non executable if any of the corre-
50 sponding resource requirements $rq_{ik}(\mu_h)$ exceeds the capacity of a renewable resource

(or the capacity of a non-renewable resource, reduced by the sum of the minimum requirements of all other activities). A non-renewable resource r_k is said to be redundant if its capacity exceeds the sum of the maximum consumption of all activities. Finally, a mode μ_h for activity a_i is inefficient if there exist some other mode μ_r with shorter duration and lower consumption for all resources. Note this is in fact a form of static dominance rule (see Section 5.2.2); as such, non-regular objectives, maximal time lags and precedence constraints other than end-to-start may prevent the rule application (see Section 2.1.3 and 4.1).

5.2.4 Multi-mode Cutset Rules

This family of rules requires to store information about past search. During the solution process, the current partial schedule is compared with the stored data; in case any solution obtainable from the current partial schedule cannot be better than the solution obtained from a previously evaluated partial schedule, then backtracking is performed. The presented formulation of the cutset rules is devised for a makespan minimization objective, but does extend to regular performance measures [76]; some care should be observed with Generalized Precedence Relations (see Sections 2.1.2 and 4.1). Cutset rules are described in [25, 37, 24, 18].

Given a partial schedule S , the cutset $C(S)$ is the set of activities scheduled so far; besides the cutset, the rule requires to store the completion time of the partial schedule $e(S)$ (i.e. the highest end time among activities in S) and the leftover capacities of all non-renewable resources. Then:

Rule 1: Dominated Heads

Let S be the partial schedule to be extended at time t in the current search step, having cutset $C(S)$; if a stored partial schedule S' exists, with: (1) the same cutset, i.e. $C(S') = C(S)$; (2) lower or equal completion time, i.e. $e(S') \leq e(S)$; (3) higher or equal leftover capacities for all non-renewable resources; then the current partial schedule needs not to be completed.

A second rule is presented in [76] and provides a bound for the time span necessary to complete the current partial schedule; the rule requires to store, for each visited partial schedule S the updated latest finish time $LET(S, a_i)$ of a_i , after all possible continuations of S have been explored.

Rule 2: Incompletable Tails

let S be the partial schedule to be extended with activity a_i at time t in the current search step, having cutset $C(S)$; if a stored partial schedule S' exists, with: (1) the same cutset, i.e. $C(S') = C(S)$; (2) higher or equal leftover capacities for all non-renewable resources; (3) $t + LET(S', a_i) - e(S') + 1 > LET(a_i)$; then the current partial schedule cannot be completed with a makespan better than the current $LET(a_i)$.

5.2.5 Effectiveness of the Bounding Rules

Some experimental evaluation of the described rules is reported in [77, 76, 37, 25]; additionally, [76] provides some details about the rule implementation. An overall thorough comparison is difficult, since different works have considered different bounding rules

1 and sometimes targeted different instance sets. As a general remark, combining bounding
 2 rules is in general fruitful, i.e. there is not sharp dominance relation.
 3

4 The non-renewable resource rule is considered to be among the most effective tech-
 5 nique and provides the highest speed-up both in [77] and [76]; the reported improvement
 6 is less substantial, but still remarkable, in [25]; this result points out the difficulties
 7 encountered by OR methods when feasible schedules are not trivial to build (e.g. when
 8 non-renewable resource capacities are taken into account). The effectiveness of the left-
 9 shift rule is also well documented; interestingly, the best results are usually obtained by
 10 testing the feasibility of local left-shifts. The single enumeration rule has a fundamental
 11 role within precedence tree branching [37]; the multi-mode rule performs nicely in the
 12 comparison from [77]. Among the cutset rules, the rule of dominated heads performs
 13 very well, definitely much better than incompletable tails. The immediate selection rule
 14 tends to be effective for small instances, but becomes more expensive and less likely to
 15 be triggered on instances with more than 10 activities. Static bounding rules are very
 16 effective for MRCPSp instances with high average resource requirements (specifically,
 17 Resource Strength – see [77,49]). The order-swap rule introduced in [37] is as effective
 18 as the local left-shift one.

20 6 Search

21 The key problem addressed by search methods for assignment and scheduling problems
 22 is how to effectively explore a search space that is in the typical case impressively
 23 large: this is a combined effect due to the domain size of temporal variables and to the
 24 presence of resource allocation choices. Efficacy can be obtained (1) by guiding search
 25 as quickly as possible towards feasible/optimal solutions; (2) by designing branching
 26 decisions so as to maximize propagation effectiveness; (3) by reducing the portion
 27 of the search space that needs to be explored (e.g. via the application of dominance
 28 rules or dynamic symmetry breaking techniques). While powerful techniques to achieve
 29 this goals are known for pure scheduling problems, it is still not-so-clear how resource
 30 assignment decisions are best treated. This section discusses the main search methods
 31 adopted in CP for resource allocation and scheduling problems (Section 6.1) and in OR
 32 for the MRCPSp (Section 6.2), trying to outline similarities and point out the main
 33 strengths/weaknesses of each approach.

36 37 6.1 Search Strategies in Constraint Programming

38 Exact CP algorithms for resource allocation and scheduling problems are usually based
 39 on Depth First Search. The nodes of the search tree represent partially instantiated
 40 schedules, where scheduled activities have fixed start time; time windows for the un-
 41 scheduled activities are maintained via the domains of start/end variables and up-
 42 dated via propagation. Search proceeds by opening choice points and posting different
 43 branching constraints along each branch; distinct strategies differ for the type of posted
 44 constraints and for the heuristics used to take non-deterministic decisions.

47 48 Two-stage search

49 Here, we refer as two-stage search to any tree search method taking into account
 50 resource assignment and scheduling variables in successive, distinct phases. This is
 51

1 the default search method for alternative resources in IBM ILOG Scheduler, where
 2 all resource assignment decisions are taken in a first stage, and a branching scheme
 3 for pure scheduling is then applied. The approach is very simple to implement, but
 4 tends to considerably under-exploit the joint propagation of precedence and resource
 5 constraints: the best results are obtained for graphs featuring a very small number of
 6 arcs. The method is also considered for the Disjunctive Temporal Problem with Finite
 7 Domain constraints [63] (see Section 3.1.2).
 8

9 *Left-Justified Random and Next Or Successor-but-not-next*

10 This search method [66] finds the smallest earliest *finish* time of all the unscheduled
 11 activities and then identifies the set of activities which are able to *start* before this time.
 12 One of the activities in this set (say a_i) is selected randomly and scheduled at its earliest
 13 start time. When backtracking, the alternative commitment is to update the earliest
 14 start time of the activity to the minimum earliest finish time of all other activities on the
 15 same resource as a_i . The approach is extended to models with alternative activities in
 16 [12], by taking into account activities for scheduling only if the corresponding execution
 17 variable ex_i is not bound to 0. When the activity is scheduled, it is simultaneously
 18 selected for execution; on backtrack, the earliest start time of the a_i is updated, but the
 19 corresponding execution variable is not modified (this ensures that search is complete).
 20

21 By scheduling activities at their earliest start time, Left Justified Random focuses
 22 on the construction of active schedules (see Section 5.2.2); this considerably narrows
 23 the search space and usually incurs no loss of optimality, since with regular performance
 24 measures the set of optimal schedules is guaranteed to contain an active one. In the
 25 experimental results in [12], Left-Justified Random reports quite poor results compared
 26 to Precedence Constraint Posting guided by texture based heuristics; we conjecture that
 27 changing the activity selection criterion may improve the outcome.

28 A closely related strategy is employed in [33] for an assignment and scheduling
 29 problem with unary resources and setup times. The method makes use of binary choice
 30 points. Let L be the set of the last activities scheduled on each resource; in the first
 31 branch, the activity a_i with minimum earliest start time is scheduled to be “next”
 32 after some activity a_j in L and assigned to the same resource. On backtracking, a_i is
 33 forced to be a “successor, but not next” and no resource assignment is performed. This
 34 strategy mainly differs from Left-Justified random for the activity selection criterion
 35 (which tends to be more effective) and for the lack of an earliest start time update on
 36 backtracking. The approach has also some analogy with Precedence Constraint Posting,
 37 since it does not require to immediately assign a start time to the selected activity.
 38

39 *Schedule Or Postpone*

40 This strategy is proposed in [59] for pure scheduling problems; at each search node an
 41 activity a_i is selected for branching; typically, this is the one with the lowest $EST(a_i)$,
 42 while $LET(a_i)$ is used to break ties. Then a binary choice point is opened and a_i is
 43 scheduled at $EST(a_i)$ along the first branch; upon backtracking the activity is marked
 44 as non-selectable (i.e. postponed) until its earliest start time is modified by propagation.
 45 In other words, the time window update performed on backtracking by Left-Justified
 46 Random is demanded here to the resource constraints: this makes the approach more
 47 general. Moreover, postponing activities reduces the size of the search space by pre-
 48 venting scheduling choices to be repeated. The main underlying idea is once again to
 49 focus search on the production of active schedules. However, according to practical
 50

51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

1 experience, the method performs well even for many non-regular objectives, despite no
 2 optimality guarantee can be given in that case.
 3

4 Intuitively, generalizing the schedule-or-postpone strategy requires to take into ac-
 5 count resource assignment decisions and to focus on the production of tight (rather
 6 than active) schedules; this is attempted in [17] within an approach with alternative
 7 resources, by adding an assignment stage before opening a schedule-or-postpone choice
 8 point. The method works very well for graphs with many precedence constraints, since
 9 in this case interleaving assignment and scheduling decisions triggers better propa-
 10 gation compared to two-stage search. Note however that, in order to produce tight
 11 schedules, it would be sufficient to postpone the activity only once *all* the possible
 12 resource assignment have been tried; as a consequence, the proposed strategy explores
 13 unnecessary portions of the search space, leaving room for future improvements.

14 *Precedence Constraint Posting (PCP)*

15 This search method has been designed for pure scheduling problems and proceeds by
 16 resolving possible resource conflicts through the addition of precedence constraints.
 17 The key idea is to identify minimal sets of activities causing a resource over-usage in
 18 case of overlapping execution, i.e. so called Minimal Critical Sets (MCS), introduced in
 19 [46] as “Minimal Forbidden Sets” and first used in CP in [22]. Due to the minimality
 20 property, the occurrence of a MCS can be prevented by adding a single precedence
 21 constraint between any two activities in the set (a so-called resolver). Branching is
 22 done either by enumerating all possible conflict resolution choices, or by opening a
 23 binary choice point where a selected resolver is alternatively posted or forbidden. MCS
 24 can be detected via (1) enumeration [53], (2) by sampling peaks over an earliest start
 25 schedule or over a worst case resource usage envelope [70], (3) by solving a minimum
 26 flow problem [61].

27 This search method has been applied to resource assignment and scheduling prob-
 28 lems in [12], modeled via alternative activities; in this case a MCS can be resolved by
 29 either posting a precedence constraint or by setting to 0 the execution variable ex_i of
 30 some activity; search is performed via binary choice points. The PCP approach allows
 31 to resolve conflicts in a non-chronological order (e.g. the hardest first); as a side effect,
 32 the resulting schedules are not necessarily active/tight (see the example in [79]) and
 33 left-shift rules are difficult to apply.

34 The actual conflict and the resolver used for branching are chosen according to some
 35 heuristic; [12] proposes in first place an adaptation of a slack based heuristic from the
 36 literature. Alternatively, the same work describes two texture based heuristics, that
 37 are in practice conflict and resolver selection procedures: the main underlying idea
 38 is to rely on some information extracted from the current problem state (a so-called
 39 texture) to identify a critical resource r_k^* and a critical time point t^* ; then sequencing
 40 decisions are taken on the activities with a chance of overlap and cause a conflict at t^* .
 41 In detail, the paper considers two textures: (1) a probabilistic resource usage profile and
 42 (2) a conflict probability function over time, from which the critical resource and time
 43 point are extracted. Both approaches perform very well in the experimental results.
 44 As a main difference with classical PCP conflict selection, the heuristics may identify
 45 critical (r_k^*, t^*) even if no conflict arises for that pair and should be therefore coupled
 46 with a conflict detection procedure. Note the worst case usage envelope employed in
 47 [70] to detect MCS can be considered a form of texture. The ideas exposed in this
 48 section are very closely related to Minimal Forbidden Set branching for the MRCPS
 49 (see Section 6.2).

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

 1 6.2 Branching Schemes for the MRCPS
 2

 3 All the exact approaches developed in an OR context for the MRCPS are based on tree
 4 search; unlike in Constraint Programming, where a search node always corresponds to
 5 a state of the domain store, branching schemes from the Operations Research literature
 6 rely on different types of schedule representation.
 7

 8 *Precedence Tree Branching*
 9

 10 This branching strategy is introduced in [81], but received a major improvement by Pat-
 11 terson in [67]. Each node of the search tree corresponds to a resource- and precedence-
 12 feasible *partial schedule*, i.e. a schedule of a set S of activities, built in chronological
 13 order from time 0. No updated information (e.g. range of possible start times) is main-
 14 tained for unscheduled activities.
 15

 16 The strategy consists in scheduling at each step of the search tree an activity
 17 whose predecessors have all been scheduled. Therefore, for each search node with par-
 18 tial schedule S a set of eligible activities $E(S)$ is univocally defined. On backtrack, a
 19 different activity is chosen, so that each path from the root of the search tree to the
 20 leaves corresponds to a possible linearization of the partial order induced on A by the
 21 precedence graph. A *mode alternative* within this branching scheme is an assignment of
 22 a mode μ_h to the target activity a_i , which is performed after the scheduling decision.
 23 On backtrack, different modes are tested, so that a scheduling decision on a_i is only
 24 undone once all modes for a_i are tested.
 25

 26 The original algorithm by Patterson is improved in [78] and [76], in particular with
 27 the introduction of bounding rules. The structure of the approach is well described in
 28 [37]. The precedence tree method suffers from symmetry issues; in particular, if two
 29 activities a_i, a_j can be independently assigned the same start time, the method will
 30 always test both enumeration sequences; this is countered by the application of the
 31 single enumeration rule (see Section 5.2), which in fact provides the highest benefits
 32 on this branching scheme.
 33

 34 *Mode and Delay Alternatives*
 35

 36 This branching method is introduced in [23, 28] for the RCPSP and is adapted to the
 37 multi-mode case in [77]; it is well described in [37] and recently used in [25]. Each node
 38 of the search tree is associated to a feasible *partial schedule* S and a time instant t .
 39 A clear distinction is then made between completed activities at time t (say $C(S, t)$)
 40 and activities in process (say $P(S, t)$); eligible activities for scheduling are those whose
 41 predecessors have all completed execution.
 42

 43 Then an attempt is made to schedule *all* eligible activities and they are added to
 44 the set of activities in process. Of course this may cause a conflict; in such a case the
 45 method branches by *withdrawing* from execution the so-called delay alternatives. Those
 46 are: (1) activities in process, i.e. in $P(S, t)$ and (2) such that, if they are removed, no
 47 resource conflict occurs. A delay alternative is called minimal if none of its proper
 48 subsets is a delay alternative; branching on minimal delay alternatives is sufficient to
 49 explore the whole search space. If no resource conflict occurs, the only minimal delay
 50 alternative is the empty set.
 51

 52 This method differs from the precedence tree based one in two regards: (1) the
 53 process branches on *sets* of activities and (2) activities scheduled at a search node may
 54 be withdrawn from execution later on: in case this is not done, the algorithm only
 55

1 builds so called non-delay schedules [79] and may miss the optimal solution (even in
 2 case the objective is regular).
 3

4 Activities are assigned a mode when they are first inserted in the $P(S, t)$ and
 5 they retain the mode assignment when they are withdrawn from execution. As a con-
 6 sequence, when the simultaneous execution of all eligible activities is probed, some
 7 activities already possess a mode, while the remaining ones are modeless. A mode al-
 8 ternative in this search strategy is a mapping that assigns a mode to every activity in
 9 the modeless eligible set; on backtracking, when the subtree corresponding to a delay
 10 alternative has been completely explored, a different mode alternative is picked.
 11

12 *Mode and Extension Alternatives*

13 This method was proposed in [80] for the RCPSP, while the adaptation to the multi-
 14 mode case is discussed in [37]. The approach is similar to mode and delay alternatives:
 15 each search node corresponds to a *partial schedule* S and a time instant t , for which
 16 sets $C(S, t)$ and $P(S, t)$ are identified. The activities with all predecessors in $C(s, t)$
 17 form the *eligible set* $E(S, t)$.
 18

19 The current partial schedule is extended by starting a *subset* of the eligible activities
 20 without violating the renewable resource constraints; conversely, in delay alternatives
 21 all eligible activities are started, then some are withdrawn from execution. In order to
 22 secure that the algorithm terminates, empty extension alternatives may be disregarded
 23 if $P(s, t) = \emptyset$ (i.e. no activity is in process). However, if there are currently activities in
 24 process, the empty set is always an extension alternative which must be tested in order
 25 to guarantee optimality; in case this is not done, the algorithm only builds non-delay
 26 schedules and may miss the optimal solution.

27 A *mode alternative* is a mapping of modes to activities and occurs as soon as
 28 they become *eligible*, before an extension alternative is selected. The backtracking
 29 mechanism is the same as for delay alternatives. This branching scheme is proven to
 30 be dominated by precedence tree and delay alternative branching in [37], at least when
 31 no bounding rule is applied.
 32

33 *Dichotomization and Time Window Tightening*

34 A branching scheme based on dichotomization is proposed for the MRCPS in [91].
 35 The approach operates on a time indexed model and is based on Special Ordered
 36 Sets (SOS, [9]); in detail, each considered SOS includes all the binary variables $E_{i,h,t}$
 37 referring to a single activity.
 38

39 Given a fractional LP solution, branching is performed by splitting the SOS re-
 40 ferring to a activity a_i , based on its average finish time t_b ; the first subset contains
 41 variables with $t \leq t_b$, the second with $t > t_b$. Two branches are defined by respec-
 42 tively setting the variables in each subset to zero. The search method is coupled with
 43 a bound-tightening step at each branch; the approach obtains interesting results.
 44

45 A related approach for the Multi-skill Project Scheduling Problem (affine to the
 46 MRCPS) is proposed by Carlier and Latapie in [68] and is based on time windows
 47 tightening; it has been adapted for Multi-skill Project Scheduling Problem by Pessan
 48 et al. [68]. At each node a time window is selected and reduced until all activities have
 49 their starting point fixed. During the process two lower bounds are used: one based
 50 on the compatibility graph that checks which activities can be scheduled at the same
 51 time, and one based on energetic reasoning that checks the mandatory parts.
 52

53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

1 ***Minimal Forbidden Sets***

2 Minimal Forbidden Sets, known in CP as Minimal Critical Sets, are introduced in [46]
3 in the early 80s; they can be used to define choice points in tree search as described
4 in Section 6.1. Unlike in the CP literature where branching is mainly done via binary
5 precedence constraints (resolvers), OR methods have explored the use of *disjunctive*
6 *precedence constraints*, i.e. by requiring a specific activity a_i in the MCS to execute
7 after *at least one* other activity a_j in the set; the specific a_j causing the delay may be
8 left undecided until all start times are assigned. The method is devised in [75] for the
9 RCPSP and is applied in [38] to the multi-mode case.

10 With this branching strategy, a search node corresponds to *fictitious* rather than to
11 a partial schedule. A fictitious schedule assigns a set of possible modes and a provisional
12 start time to each activity of the project; the resulting representation is very similar
13 to the one obtained in CP via the domain store. Disjunctive precedence constraints
14 are based on the same idea of delay alternatives, but enable one to consider conflicts
15 in non-chronological order; in particular, the method described in [38] systematically
16 branches on the (estimated) hardest decision. The specific rule applied depends on the
17 type of this decision; if this is a mode decision, each branch is a mode assignment for the
18 corresponding activity. If it is a conflict decision, each branch is a possible disjunctive
19 precedence constraint to resolve the conflict.

20 The use of disjunctive precedence relations allows a remarkable reduction of the
21 search tree compared to binary resolvers; as a drawback, with this type of precedence
22 relations the feasible space becomes a *union* of convex polyhedra and Generalized Arc
23 Consistency on the temporal constraints can no longer be achieved in polynomial time.

25 **7 Decomposition based Solution Methods**

26 Most of the decomposition approaches for assignment and scheduling problems are
27 strongly related to Logic based Benders' Decomposition (LBD), introduced in [45]
28 and formalized in [39, 42]. LBD generalizes classic Benders' Decomposition [15] and
29 considerably broadens its application field, by replacing the linear dual used for cut
30 generation with a so-called *inference dual*. The method breaks a combinatorial problem
31 into a master and a subproblem, which are solved in sequence; the master solution is
32 used to prime the subproblem, then a cut is generated (Benders' cut) and permanently
33 added to the master problem; the process repeats until the master and the sub-problem
34 converge in value and optimality is proved. In the context of allocation and scheduling,
35 resource assignment variables x_{ik} are typically included in the master problem, while
36 scheduling variables s_i, e_i only appear in the subproblem; assuming a minimization
37 objective, the cost function $F(\bar{x}, \bar{s}, \bar{e})$ is replaced by a relaxed function $G(\bar{x})$ such that
38 $G(\bar{x}) \leq F(\bar{x}, \bar{s}, \bar{e})$. The approach allows one to use heterogeneous solution techniques:
39 quite often MILP is employed with the master and CP with the subproblem, due to
40 the effectiveness of CP constraints and search strategies for pure scheduling.

41 Benders' cuts at each iteration, once an assignment of master problem variables \bar{x}' is
42 available, are obtained from the solution of an inference dual. In detail, let the function
43 $F^*(\bar{x})$ denote the minimum value of $F(\bar{x}, \bar{s}, \bar{e})$ corresponding to an assignment of the
44 master problem variables: the inference dual consists in computing the value $F^*(\bar{x}')$.
45 Then, the algorithm designer is responsible to identify a bounding function $\beta_{\bar{x}'}(\bar{x})$ such
46 that $\beta_{\bar{x}'}(\bar{x}) \leq F^*(\bar{x})$ and $\beta_{\bar{x}'}(\bar{x}') = F^*(\bar{x}')$; in other words, $\beta_{\bar{x}'}(\bar{x})$ equals the optimal
47 $F(\bar{x}, \bar{s}, \bar{e})$ when $\bar{x} = \bar{x}'$, otherwise the function provides a valid bound. The subscript
48

1 of the bounding function denotes the \bar{x} values used for its construction. The Benders'
 2 cut consists in forcing the master-problem objective to be smaller than the bounding
 3 function, i.e. $F(\bar{x}) < \beta_{\bar{x}'}(\bar{x})$. At each iteration, the cumulated Benders' cuts prevent
 4 old master problem solutions from being re-proposed.
 5

6 The approach was first applied to allocation and scheduling problems in [36, 47],
 7 obtaining dramatic improvements over pure CP/MILP methods and over a hybrid algo-
 8 rithm using CP for branching and LP-based lower bounds. Other application examples
 9 include [43] for tardiness costs, [21] for scheduling on hard-real time systems, [72] for
 10 power consumption minimization on multi-core CPUs.

11 In case of resource based objectives (see Section 2.1.3) the sub-problem has no cost
 12 function, making CP an even better candidate for its solution; this is well documented
 13 in [47, 44]. In that case, Benders' cuts can be simply formulated as constraints on the
 14 x_{ik} variables rather than bounds on the master cost function. Sometimes, it is possible
 15 to further decouple the subproblem into independent components (e.g. a scheduling
 16 problem for each resource), with strong efficiency improvements: this typically occurs
 17 when no precedence constraint is specified and is discussed in Section 3.3. In general, a
 18 tighter bounding function $\beta_{\bar{x}'}(\bar{x})$ results in quicker convergence, providing motivation
 19 for cut strengthening techniques: an efficient ad hoc method is reported in [40].
 20

21 *Balancing*

22 Quite often, the master problem turns out to be considerably heavier than the sub-
 23 problem, therefore limiting the number of full iterations performed in a given amount
 24 of time. This is a relevant issue and can be dealt with in a number of ways: (1) by not
 25 solving the master problem to optimality: this especially makes sense in case of schedul-
 26 ing dependent objectives, where optimal values for the *relaxed* cost function $F(\bar{x})$ do
 27 not necessarily result in better solutions; this is observed in [82] and provides some
 28 arguments for the use of CP in the master problem. (2) By solving the sub-problem
 29 while still branching to search for a master solution: this is key idea behind so-called
 30 branch-and-check [82] and requires a method to formulate the sub-problem when some
 31 master problem variables are unassigned – e.g. the master and the sub-problem may
 32 share some variables. (3) By using stronger and more computation expensive cuts: for
 33 example, improved cuts can be obtained via explanation minimization, by repeatedly
 34 solving polynomial complexity sub-problems [21] or even NP-hard ones [62]. (4) By
 35 applying LBD recursively on the master problem itself, obtaining a multi-stage de-
 36 composition: [17] describes the approach in a specific application settings and provides
 37 empirical rules to balance the computational load of the different decomposition stages.
 38

39 *Subproblem Relaxation*

40 Logic based Benders' Decomposition incurs the risk of loosing valuable information due
 41 to decoupling; addressing this issue is the fundamental problem during the design of a
 42 LBD approach: this is of course tackled by devising effective cuts, but the interaction
 43 between the two stages can be further tightened by introducing in the master problem a
 44 so-called *subproblem relaxation*. This may consist of additional constraints (i.e. some
 45 of the subproblem constraints, possibly relaxed) or may be an additional bounding
 46 function on the master problem objective. The actual formulation has to be given case
 47 by case: examples can be found in basically all LBD related references in this paper.
 48 The use of an effective subproblem relaxation often has a tremendous impact on the
 49 performance: in [82] it is shown how removing the relaxation from the model in [47]
 50 increases the search time by several orders of magnitude.
 51

52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 8 Conclusions

We provided an overview of state of the art approaches for a class of resource allocation and scheduling problems, arising in a variety of real world settings. Given the amount of problem variants we chose to focus this work on techniques to address individual problem traits, rather than on devising an exhaustive (most likely too complex) classification. In particular, we mainly drew the presented pool of algorithms and methods from scheduling related OR and CP literature. Constraint Programming is a natural candidate to support the integration of heterogeneous techniques: its typical distinction between model, propagation and search provided the backbone for the work organization.

Hybrid methods were given prominent importance, as they proved to be particularly effective on allocation and scheduling problems. On the CP side, we observed a lack of effective algorithms for joint filtering on assignment and scheduling variables: we find this stimulating rather than disappointing, as it may provide a nice benchmark for novel propagation techniques, based on decomposition, SAT hybridization or alternative representations of the domain store. Indeed, the whole paper can be seen as an attempt to provide hints for the development of new hybrid algorithms and novel filtering methods.

We decided to limit our discussion to exact approaches; however, given the impressive complexity of this class of problems, a very large number of works from the literature adopts heuristic solution methods. We hope our effort will provide motivation to the research community for surveying the state of the art of heuristics for resource assignment and scheduling.

References

1. P. Baptiste, P. Laborie, C. Le Pape, and W. Nuijten. Constraint-based scheduling and planning. *Foundations of Artificial Intelligence*, 2:761–799, 2006.
2. P. Baptiste and C. Le Pape. Disjunctive constraints for manufacturing scheduling: Principles and extensions. *International Journal of Computer Integrated Manufacturing*, 9(4):306–310, 1996.
3. P. Baptiste, C. Le Pape, and W. Nuijten. *Constraint-based scheduling*. Kluwer Academic Publishers, 2001.
4. R. Barták and O. Čepek. Nested Precedence Networks with Alternatives: Recognition, Tractability, and Models. In *Proc. of AIMSA*, pages 235–246, 2008.
5. R. Barták, O. Čepek, and P. Surynek. Modelling Alternatives in Temporal Networks. In *Proc. of IEEE SCIS*, pages 129–136, 2007.
6. R. Barták, O. Čepek, and M. Hejna. Temporal Reasoning in Nested Temporal Networks with Alternatives. *Recent Advances in Constraints*, pages 17–31, 2008.
7. R. Barták, O. Čepek, and P. Surynek. Discovering implied constraints in precedence graphs with alternatives. *Annals of Operations Research*, 180(1):233–263, 2008.
8. M. Bartusch, R. H. Möhring, and F. J. Radermacher. Scheduling project networks with resource constraints and time windows. *Annals of Operations Research*, 16(1):199–240, 1988.
9. E. M. L. Beale and J. J. H. Forrest. Global optimization using special ordered sets. *Mathematical Programming*, 10(1):52–69, 1976.
10. JE Beasley and M Krishnamoorthy. Scheduling aircraft landings-the static case. *Transportation science*, 34(2):180–197, 2000.
11. J. C. Beck and M. S. Fox. Scheduling Alternative Activities. In *Proc. of AAAI/IAAI*, pages 680–687, 1999.
12. J. C. Beck and M. S. Fox. Constraint-directed techniques for scheduling alternative activities. *Artificial Intelligence*, 121(1-2):211–250, 2000.

- 1 13. Nicolas Beldiceanu, Mats Carlsson, Sophie Demassey, and Emmanuel Poder. New filtering
2 for the cumulative constraint in the context of non-overlapping rectangles. *Annals of*
3 *Operations Research*, 184(1):27–50, March 2011.
- 4 14. Odile Bellenguez-Morineau and Emmanuel Néron. A Branch-and-Bound method for solv-
5 ing Multi-Skill Project Scheduling Problem. *RAIRO - Operations Research*, 41(02):155–
6 170, April 2007.
- 7 15. J. F. Benders. Partitioning procedures for solving mixed-variables programming problems.
8 *Numerische Mathematik*, 4(1):238–252, 1962.
- 9 16. L. Benini, D. Bertozzi, A. Guerri, and M. Milano. Allocation and scheduling for MPSoCs
10 via decomposition and no-good generation. *Proc. of CP*, pages 107–121, 2005.
- 11 17. L. Benini, M. Lombardi, M. Milano, and M. Ruggiero. Optimal Allocation and Scheduling
12 for the Cell BE Platform. *Annals of Operations Research*, 184(1):51–77, 2011.
- 13 18. P. Brucker, A. Drexel, R. H. Möhring, K. Neumann, and E. Pesch. Resource-constrained
14 project scheduling: Notation, classification, models, and methods. *European Journal of*
15 *Operational Research*, 112(1):3–41, 1999.
- 16 19. J. Buddhakulsomsiri and D. Kim. Properties of multi-mode resource-constrained project
17 scheduling problems with resource vacations and activity splitting. *European Journal of*
18 *Operational Research*, 175(1):279–295, November 2006.
- 19 20. R. E. Burkard, M. Dell’Amico, and S. Martello. *Assignment problems*. Society for Indus-
20 trial Mathematics, 2009.
- 21 21. H. Cambazard, P.E. Hladik, A.M. Déplanche, N. Jussien, and Y. Trinquet. Decomposition
22 and learning for a hard real time task allocation problem. *Proc. of CP*, pages 153–167,
23 2004.
- 24 22. Amedeo Cesta, Angelo Oddi, and S.F. Smith. Scheduling Multi-Capacitated Resources
25 under Complex Temporal Constraints. *Proc. of CP*, pages 465–465, 1998.
- 26 23. N. Christofides, R. Alvarez-Valdes, and J. M. Tamarit. Project scheduling with resource
27 constraints: A branch and bound approach. *European Journal of Operational Research*,
28 29(3):262–273, 1987.
- 29 24. B. De Reyck, E. Demeulemeester, and W. Herroelen. Local search methods for the discrete
30 time/resource trade-off problem in project networks. *Naval Research Logistics*, 45(6):553–
31 578, 1998.
- 32 25. F. Deblaere, E. Demeulemeester, and W. Herroelen. Reactive scheduling in the multi-mode
33 RCPSP. *Computers & Operations Research*, pages 1–12, January 2010.
- 34 26. R. Dechter, I. Meiri, and J. Pearl. Temporal Constraint Networks. *Artificial Intelligence*,
35 49(1-3):61–95, 1991.
- 36 27. E. Demeulemeester, B. De Reyck, and W. Herroelen. The discrete time/resource
37 trade-off problem in project networks: a branch-and-bound approach. *IIE transactions*,
38 32(11):1059–1069, 2000.
- 39 28. E. L. Demeulemeester and W. Herroelen. A branch-and-bound procedure for multiple
40 resource-constrained project scheduling problem. *Management science*, 38(12):1803–1818,
41 1992.
- 42 29. A. Drexel and J. Gruenewald. Nonpreemptive multi-mode resource-constrained project
43 scheduling. *IIE transactions*, 25(5):74–81, 1993.
- 44 30. S. E. Elmaghraby and J. Kamburowski. The Analysis of Activity Networks Under Gener-
45 alized Precedence Relations (GPRs). *Management Science*, 38(9):1245–1263, 1992.
- 46 31. S.E. Elmaghraby. *Activity networks: Project planning and control by network models*.
47 Wiley New York, 1977.
- 48 32. A.T. Ernst, M. Krishnamoorthy, and R.H. Storer. Heuristic and exact algorithms for
49 scheduling aircraft landings. *Networks*, 34(3):229–241, 1999.
- 50 33. F. Focacci, P. Laborie, and W. Nuijten. Solving scheduling problems with setup times and
51 alternative resources. In *Proc. of ICAPS*, 2000.
- 52 34. F. Focacci, A. Lodi, and M. Milano. Cost-based domain filtering. In *Proc. of CP*, pages
53 189–203, 1999.
- 54 35. LM Gambardella and M. Mastrolilli. Effective neighborhood functions for the flexible job
55 shop problem. *Journal of Scheduling*, 3:3, 1996.
- 56 36. I. Harjunkoski, V. Jain, and I.E. Grossman. Hybrid mixed-integer/constraint logic pro-
57 gramming strategies for solving scheduling and combinatorial optimization problems.
58 *Computers & Chemical Engineering*, 24(2-7):337–343, 2000.
- 59 37. S. Hartmann and A. Drexel. Project scheduling with multiple modes: A comparison of
60 exact algorithms. *Networks*, 32(4):283–297, 1998.
- 61
- 62
- 63
- 64
- 65

-
- 1 38. R. Heilmann. A branch-and-bound procedure for the multi-mode resource-constrained
2 project scheduling problem with minimum and maximum time lags. *European Journal of*
3 *Operational Research*, 144(2):348–365, January 2003.
4 39. J. Hooker. *Logic-based methods for optimization: combining optimization and constraint*
5 *satisfaction*. John Wiley & Sons, 2000.
6 40. J. N. Hooker. A Hybrid Method for Planning and Scheduling. *Constraints*, 10(4):385–401,
7 2005.
8 41. J. N. Hooker. An Integrated Method for Planning and Scheduling to Minimize Tardiness.
9 *Constraints*, 11(2-3):139–157, 2006.
10 42. J. N. Hooker and G. Ottosson. Logic-based Benders decomposition. *Mathematical Pro-*
11 *gramming*, 96(1):33–60, 2003.
12 43. J.N. Hooker. Planning and scheduling to minimize tardiness. *Proc. of CP*, 3709:314–327,
13 2005.
14 44. J.N. Hooker. Planning and scheduling by logic-based benders decomposition. *Operations*
15 *Research*, 55(3):588, 2007.
16 45. John N Hooker and H Yan. Verifying logic circuits by Benders decomposition. *Proc. of*
17 *CP*, pages 267–288, 1995.
18 46. G. Igelmund and F. J. Radermacher. Preselective strategies for the optimization of stochas-
19 tic project networks under resource constraints. *Networks*, 13(1):1–28, January 1983.
20 47. V. Jain and I. E. Grossmann. Algorithms for Hybrid MILP / CP Models for a Class of
21 Optimization Problems. *INFORMS Journal on Computing*, pages 258–276, 2001.
22 48. J.E. Kelley and M. R. Walker. Critical-path planning and scheduling. In *Proc. of eastern*
23 *joint IRE-AIEE-ACM conference*, pages 160–173, New York, USA, 1959. ACM.
24 49. R. Kolisch. PSPLIB - A project scheduling problem library. *European Journal of Opera-*
25 *tional Research*, 96(1):205–216, January 1997.
26 50. R. Kolisch and A. Drexl. Local search for nonpreemptive multi-mode resource-constrained
27 project scheduling. *IIE Transactions*, 29(11):987–999, November 1997.
28 51. Y. K. Kwok. Dynamic critical-path scheduling: An effective technique for allocating task
29 graphs to multiprocessors. *Parallel and Distributed Systems, IEEE*, 7(5):506—521, 1996.
30 52. Y. K. Kwok and I. Ahmad. Static scheduling algorithms for allocating directed task graphs
31 to multiprocessors. *ACM Computing Surveys (CSUR)*, 31(4):406–471, 1999.
32 53. P. Laborie. Complete MCS-Based Search: Application to Resource Constrained Project
33 Scheduling. In *Proc. of IJCAI*, pages 181–186. Professional Book Center, 2005.
34 54. P. Laborie. IBM ILOG CP Optimizer for detailed scheduling illustrated on three problems.
35 In *Proc. of CPAIOR*, pages 148–162, 2009.
36 55. P. Laborie and D. Godard. Self-adapting large neighborhood search: Application to single-
37 mode scheduling problems. In *Proc. of MISTA*, 2007.
38 56. P. Laborie and J. Rogerie. Reasoning with conditional time-intervals. In *Proc. of FLAIRS*,
39 pages 555–560, 2008.
40 57. P. Laborie, J. Rogerie, P. Shaw, and P. Vilim. Reasoning with Conditional Time-Intervals
41 Part II: An Algebraical Model for Resources. In *Proc. of FLAIRS*, pages 201–206, 2009.
42 58. C. Le Pape. Using a constraint-based scheduling library to solve a specific scheduling
43 problem. In *Proc. of AAAI-SIGMAN*, 1994.
44 59. C. Le Pape, P. Couronné, D. Vergamini, and V. Gosselin. Time-versus-capacity compro-
45 mises in project scheduling. *AISB QUARTERLY*, page 19, 1995.
46 60. R Leupers. Instruction scheduling for clustered VLIW DSPs. *Proc. of PACT*, pages 291–,
47 2000.
48 61. M. Lombardi and M. Milano. A Precedence Constraint Posting Approach for the RCPSP
49 with Time Lags and Variable Durations. In *Proc. of CP*, pages 569–583, 2009.
50 62. M. Lombardi and M. Milano. Allocation and scheduling of Conditional Task Graphs.
51 *Artificial Intelligence*, 174(7-8):500–529, 2010.
52 63. M. D. Moffitt, B. Peintner, and M. E. Pollack. Augmenting disjunctive temporal problems
53 with finite-domain constraints. In *Proc. of AAAI*, pages 1–6, 2005.
54 64. W. Nuijten. *Time and resource constrained scheduling : a constraint satisfaction approach*.
55 PhD thesis, Technische Universiteit Eindhoven, 1994.
56 65. W. Nuijten, T. Bousonville, F. Focacci, D. Godard, and C. Le Pape. Towards an industrial
57 manufacturing scheduling problem and test bed. In *Proc. of PMS*, 2004.
58 66. W. P. M. Nuijten, E. H. L. Aarts, D.A.A. van Arp Talmaan Kip, and K.M. van Hee.
59 Randomized constraint satisfaction for job shop scheduling. In *Proc. of IJCAI*, pages
60 251–262, Chambery, France, 1993.

-
- 1 67. J.H. Patterson, R. Slowinski, F.B. Talbot, and J. Weglarz. An algorithm for a general
2 class of precedence and resource constrained scheduling problems. *Advances in project*
3 *scheduling*, pages 3–28, 1989.
4 68. C. Pesson, O. Bellenguez-Morineau, and E. Néron. Multi-skill Project Scheduling Problem
5 and Total Productive Maintenance. *Proceedings of MISTA*, pages 608–610, 2007.
6 69. E. Poder and N. Beldiceanu. Filtering for a Continuous Multi-Resources cumulative Con-
7 straint with Resource Consumption and Production. In *Proc. of ICAPS*, pages 264–271,
8 2008.
9 70. N. Policella, A. Cesta, A. Oddi, and S. F. Smith. From precedence constraint posting
10 to partial order schedules: A CSP approach to Robust Scheduling. *AI Communications*,
11 20(3):163–180, 2007.
12 71. F. J. Radermacher. Scheduling of project networks. *Annals of Operations Research*,
13 4(1):227–252, 1985.
14 72. M. Ruggiero, P. Gioia, G. Alessio, L. Benini, M. Michela, Davide Bertozzi, and Alexandru
15 Andrei. A cooperative, accurate solving framework for optimal allocation, scheduling and
16 frequency selection on energy-efficient mpsocs. In *Proc. of SoC*, pages 1–4. IEEE, 2007.
17 73. M. Sabzehtparvar and S. M. Seyed-Hosseini. A mathematical model for the multi-mode
18 resource-constrained project scheduling problem with mode dependent time lags. *The*
19 *Journal of Supercomputing*, 44(3):257–273, November 2007.
20 74. A. Schutt, T. Feydy, P. J. Stuckey, and M. Wallace. Why Cumulative Decomposition Is
21 Not as Bad as It Sounds. In *Proc. of CP*, pages 746–761, 2009.
22 75. C. Schwindt. *Verfahren zur Lösung des ressourcenbeschränkten Projektdauermin-
23 imierungsproblems mit planungsbabhängigen Zeitfenstern*. Shaker, 1998.
24 76. A. Sprecher and A. Drexl. Multi-mode resource-constrained project scheduling by a simple,
25 general and powerful sequencing algorithm. *European Journal of Operational Research*,
26 107(2):431–450, 1998.
27 77. A. Sprecher, S. Hartmann, and A. Drexl. An exact algorithm for project scheduling with
28 multiple modes. *OR Spectrum*, 19(3):195–203, 1997.
29 78. A. Sprecher and C. L. Hwang. *Resource-constrained project scheduling: exact methods for
30 the multi-mode case*. Springer, 1994.
31 79. A. Sprecher, R. Kolisch, and A. Drexl. Semi-active, active, and non-delay schedules for
32 the resource-constrained project scheduling problem. *European Journal of Operational
33 Research*, 80(1):94–102, 1995.
34 80. J. P. Stinson, E. W. Davis, and B. M. Khumawala. Multiple Resource-Constrained Schedul-
35 ing Using Branch and Bound. *IIE Transactions*, 10(3):252–259, 1978.
36 81. F. B. Talbot. Resource-constrained project scheduling with time-resource tradeoffs: The
37 nonpreemptive case. *Management Science*, 28(10):1197–1210, 1982.
38 82. E S Thorsteinsson. Branch and Check: A hybrid framework integrating mixed integer
39 programming and constraint programming. In *Proc. of CP*, pages 16–30, Paphos, Cyprus,
40 November 2001.
41 83. C. Timpe. Solving planning and scheduling problems with combined integer and constraint
42 programming. *OR spectrum*, 24(4):431–448, 2002.
43 84. P. Toth and D. Vigo. *The vehicle routing problem*, volume 9. Society for Industrial
44 Mathematics, 2002.
45 85. V. Van Peteghem and M. Vanhoucke. A genetic algorithm for the preemptive and non-
46 preemptive multi-mode resource-constrained project scheduling problem. *European Jour-
47 nal of Operational Research*, 201(2):409–418, March 2010.
48 86. P. Vilím. Max energy filtering algorithm for discrete cumulative resources. *Proc. of
49 CPAIOR*, pages 294–308, 2009.
50 87. P. Vilím. Timetable Edge Finding Filtering Algorithm for Discrete Cumulative Resources.
51 *Proc. of CPAIOR*, pages 230–245, 2011.
52 88. P. Vilím, R. Barták, and O. Cepel. Extension of (n log n) Filtering Algorithms for the
53 Unary Resource Constraint to Optional Activities. *Constraints*, 10(4):403–425, 2005.
54 89. S. Vo and A. Witt. Hybrid flow shop scheduling as a multi-mode multi-project scheduling
55 problem with batching requirements: a real-world application. *International Journal of
56 Production Economics*, 105(2), 2007.
57 90. J. C. Zapata, B. M. Hodge, and G. V. Reklaitis. The Multimode Resource Constrained
58 Multiproject Scheduling Problem : Alternative Formulations. *AIChE Journal*, 54(8), 2008.
59 91. G. Zhu, J. F. Bard, and G. Yu. A Branch-and-Cut Procedure for the Multimode Resource-
60 Constrained Project-Scheduling Problem. *INFORMS Journal on Computing*, 18(3):377–
61 390, January 2006.