



## Discrete Optimization

## A combinatorial flow-based formulation for temporal bin packing problems

J. Martinovic<sup>a,\*</sup>, N. Strasdat<sup>a</sup>, J. Valério de Carvalho<sup>b</sup>, F. Furini<sup>c</sup><sup>a</sup>Institute of Numerical Mathematics, Technische Universität Dresden, 01062 Dresden, Germany<sup>b</sup>Departamento de Produção e Sistemas/Centro ALGORITMI, Universidade do Minho, 4710-057 Braga, Portugal<sup>c</sup>Department of Computer, Control and Management Engineering “Antonio Ruberti”, Sapienza University of Rome, 00185 Roma, Italy

## ARTICLE INFO

## Article history:

Received 23 May 2022

Accepted 5 October 2022

Available online 13 October 2022

## Keywords:

Combinatorial optimization

Temporal bin packing

Fire ups

Interval scheduling

Flow formulation

## ABSTRACT

We consider two neighboring generalizations of the classical bin packing problem: the temporal bin packing problem (TBPP) and the temporal bin packing problem with fire-ups (TBPP-FU). In both cases, the task is to arrange a set of given jobs, characterized by a resource consumption and an activity window, on homogeneous servers of limited capacity. To keep operational costs but also energy consumption low, TBPP is concerned with minimizing the number of servers in use, whereas TBPP-FU additionally takes into account the switch-on processes required for their operation. Either way, challenging integer optimization problems are obtained, which can differ significantly from each other despite the seemingly only marginal variation of the problems. In the literature, a branch-and-price method enriched with many preprocessing steps (for TBPP) and compact formulations (for TBPP-FU), benefiting from numerous reduction methods, have emerged as, currently, the most promising solution methods. In this paper, we introduce, in a sense, a unified solution framework for both problems (and, in fact, a wide variety of further interval scheduling applications) based on graph theory. Any scientific contributions in this direction failed so far because of the exponential size of the associated networks. The approach we present in this article does not change the theoretical exponentiality itself, but it can make it controllable by clever construction of the resulting graphs. In particular, for the first time all classical benchmark instances (and even larger ones) for the two problems can be solved – in times that significantly improve those of the previous approaches.

© 2022 Elsevier B.V. All rights reserved.

## 1. Introduction

## 1.1. General Overview

The optimal assignment of given jobs to one or more servers with limited capacity is an important theoretical problem in discrete optimization, but also highly relevant in many applications from computer science (Bartlett et al., 2005), logistics (Hall & Magazine, 1994), or communications engineering (Chen, Hassin, & Tzur, 2002). Despite some clear relationships between the underlying abstract problems, a wide variety of different specifications and associated terminologies have developed independently in recent years, in each of these scientific fields. To provide a coherent overview, in this article we would like to mainly focus on the operations research (or rather, the cutting and packing) perspective,

but we will also refer to important concepts and results from the other areas mentioned above for further information.

In most of these scheduling problems, we consider a set of  $n \in \mathbb{N}$  given jobs, each characterized by a profit  $p_i$ , a resource consumption  $c_i$ , and an activity interval  $[s_i, e_i)$  with  $s_i < e_i$ , that have to be assigned to a single- or multi-server architecture the capacity  $C$  of which has to be respected at all instants of time. Note that, given these assumptions, jobs cannot move in time, so that the scheduling is only done with respect to the capacity dimension (that is, the machines). Without loss of generality, all input data are assumed to be nonnegative integer numbers. Classical decision-making problems arising in this context have been coarsely classified in Angelelli, Bianchessi, & Filippi (2014) as follows:

- Q1: Is it possible to arrange all jobs on a fixed set of  $r \in \mathbb{N}$  given servers?  
 Q2: Which is the subset of jobs that yields the largest profit when arranged on  $r \in \mathbb{N}$  given servers?  
 Q3: Which is the smallest number  $r \in \mathbb{N}$  of servers needed to arrange all jobs?

\* Corresponding author.

E-mail addresses: [john.martinovic@tu-dresden.de](mailto:john.martinovic@tu-dresden.de) (J. Martinovic), [nico.strasdat@tu-dresden.de](mailto:nico.strasdat@tu-dresden.de) (N. Strasdat), [vc@dps.uminho.pt](mailto:vc@dps.uminho.pt) (J. Valério de Carvalho), [fabio.furini@uniroma1.it](mailto:fabio.furini@uniroma1.it) (F. Furini).

Although Q1 can be interpreted as a decision version of Q3, and thus there are strong relationships between these two questions, the relevant literature has initially focussed intensively on Q2. Probably the most significant special case of that question is answered by the *temporal knapsack problem* (TKP), where a profit-maximal allocation of a single server (i.e.,  $r = 1$ ) has to be found. However, referring to concrete applications in communications theory, this problem was originally introduced as the *bandwidth allocation problem* (BAP), see Bar-Noy, Canetti, Kutten, Mansour, & Schieber (1999), Chen et al. (2002), generalizing some preliminary concepts from an earlier publication (Arkin & Silverberg, 1987) in machine scheduling. In that framework, a server can also be thought of as a fixed-capacity communication channel for which there are certain requests to reserve bandwidth (that is, to transmit information). Consequently, assuming the profit of each job to be proportional to its area in the capacity-time space, its contribution to the objective function is directly linked to the amount of information conveyed by it, see Bar-Noy et al. (1999). Given the state of computational hardware and commercial software (for the exact solution of such problems) at that time, the literature initially focussed on heuristic methods and corresponding performance analyses, but also on complexity-theoretic aspects of the overall problem or of special cases. For the latter, important milestones can be summarized as follows:

- The TKP with *uniform weights*, i.e., a scenario where all  $c_i$  are identical, was shown to be polynomially solvable in  $\mathcal{O}(n^2 \log n)$  in Arkin & Silverberg (1987, Theorem 1) by drawing connections to the coloring of interval graphs and minimum cost flow problems, respectively<sup>1</sup>. In that special case, the capacity constraint reduces to a cardinality constraint, and the problem under consideration is also referred to as the *interval scheduling problem*, see Kolen, Lenstra, Papadimitriou, & Spiessma (2007) for a good survey article.
- In contrast, the TKP with *uniform profits*, where  $p_i$  is assumed to be identical for all jobs, turns out to be  $\mathcal{NP}$ -hard, see Darmann, Pferschy, & Schauer (2010, Theorem 1) for a proof drawing a connection to a specific partition problem. For the sake of completeness, we mention that the TKP is called the *resource allocation problem* (RAP) in that publication.
- In Chen et al. (2002, Theorem 1), the TKP was shown to be polynomially solvable if the capacity  $C$  is not part of the input. To this end, an  $\mathcal{O}(n^{C+1})$  time algorithm based on dynamic programming was described, establishing some early foundations for a graph-theoretic interpretation of the TKP. Note that a similar observation was already part of Arkin & Silverberg (1987, Theorem 3), but this result did not refer to the traditional TKP.
- For the general TKP, significant contributions were made in Chen et al. (2002). Here, the authors distinguish between concrete specifications in bandwidth or memory allocation. The main difference is that in the second application, the so-called *storage allocation problem* (SAP), the jobs must be positioned as actual rectangles, i.e., they particularly must consume contiguous capacity units at any time. This represents a fundamental difference to the BAP, but both application examples are reminiscent of classical two-dimensional assignment problems of cutting and packing. In fact, the authors also point out strong relations to the *multidimensional knapsack problem*, see Kellerer, Pferschy, & Pisinger (2004) for a very good overview. However, the dimension (and, thus, also the size) of this substitute problem highly depends on the input data of the TKP (more pre-

cisely, the time horizon), so that both problems (BAP and SAP) turn out to be  $\mathcal{NP}$ -hard even in the very restricted case when  $s_i = 0$ ,  $e_i = 1$ , and  $p_i = c_i$  hold for all jobs.

As can be seen from this overview, already the TKP (which did not yet appear under this name in the aforementioned publications) is a very challenging problem from a theoretical point of view. For this reason, heuristic methods have been an essential subject of further investigations. We refer the interested reader to some classic approaches proposed in Bar-Noy, Bar-Yehuda, Freund, & Naor (2001), Chen et al. (2002), Darmann et al. (2010) and to algorithms for handling more general problem aspects such as on-line scenarios with preemption (Bar-Noy et al., 1999) and specific job durations that can be smaller than the activity interval, so that there is more flexibility in allocation (Phillips, Uma, & Wein, 2000).

## 1.2. Problem-specific Literature Review

From a more mathematically-oriented perspective, the TKP was formally introduced in an article addressing an application in the context of resource allocation in high performance computing, see Bartlett et al. (2005). In that work, the authors suggested tackling the TKP by techniques combining constraint programming and branch-and-cut, but their algorithms were not as strong as a direct solution of the integer program by the commercial CPLEX-solver. Some years later, however, solution methods for the TKP were significantly advanced by applying Dantzig-Wolfe decomposition methods (Caprara, Furini, & Malaguti, 2013; Caprara, Furini, Malaguti, & Traversi, 2016). Remarkably, in that strategy, the (relatively large) set of constraints is first systematically partitioned, then each of these (relatively small) partition classes is separately convexified, and finally the partial solutions obtained in this way are harmonized. In that regard, it is of particular benefit that an arbitrary but fixed variable occurs exclusively in successive constraints, see also (Furini, 2011, Chapter 3). Meanwhile, further refinements of such decomposition methods have been discussed in the literature, see Clausen, Lusby, & Ropke (2022), Gschwind & Irnich (2017) for two recent examples. Moreover, a standard Dantzig-Wolfe decomposition also empirically proved to be very helpful for a multi-server version of TKP, referred to as the *operational interval scheduling with a resource constraint* (ORSIC) in Angelelli et al. (2014). However, as the authors admit, even this approach does not address the question Q3 from the above list, which is identified as interesting future work in the concluding section of Angelelli et al. (2014), but has remained untouched (at least in terms of the efficiency of the solution approaches) for quite a long time in the scientific community.

To this end, the current paper focuses on the optimization problem pertinent to that open question, namely the *temporal bin packing problem* (TBPP), which consists of finding the minimum number of servers required to accommodate all given jobs. Even though it is a rather obvious generalization of the well-known *bin packing problem* (BPP), see Delorme, Iori, & Martello (2016), Scheithauer (2018), Valério de Carvalho (2002) for some very good and thorough overviews, the TBPP has recently been mentioned for the first time in the context of a concrete practical application from computer science in de Cauwer, Mehta, & O'Sullivan (2016). In fact, the alarming predictions concerning the exponential increase in the energy consumption of physical computing resources, illustrated in recent studies like (Barnett, Jain, Andra, & Khurana, 2018), have made industry and scientific communities take notice, and consequently prompted a sustained intensification of the debate and research on more energy-efficient operating options, see Buyya, Yeo, Venugopal, Broberg, & Brandic (2009), Kaplan, Forrest, & Kindler (2008) for some general aspects and Fettweis et al. (2019) for some specific ideas and results of a leading European research cluster.

<sup>1</sup> In fact, Arkin & Silverberg (1987) consider a slightly different problem description with  $C$  identical machines having unit capacity, and  $c_i = 1$  for all  $i \in I$ . However, on closer examination, this is nothing else than a TKP on a single server with capacity  $C$ .

Given its relations to other already well-studied optimization problems (like strip packing), which exist but are less helpful in detail, see [Martinovic, Strasdat, & Selch \(2021, Section 1\)](#), the consideration of exact solution methods for the TBPP represents an independent branch of research, see [Dell'Amico, Furini, & Iori \(2020\)](#) for a first rigorous investigation of compact models and further more sophisticated methods. As a result of that research, currently, the most promising algorithm for solving the TBPP is a branch-and-price method that incorporates numerous lower bounds and heuristics and leads to convincing results in numerical test calculations. Despite all these efforts, still not all associated benchmark instances can be solved optimally in reasonable time, as reported in [Dell'Amico et al. \(2020\)](#).

The last statement also applies, and even more clearly, to the much younger *temporal bin packing problem with fire-ups (TBPP-FU)*. In that scenario, we assume the same input data as for the TBPP, but in the objective function we have to minimize a weighted sum of the number of servers in use and the number of switch-on processes (so-called *fire-ups*) required during operation. More precisely, a fire-up has to be counted whenever a server goes from an inactive state (no load at all) into an active state (executing some job). Including this additional aspect generally leads to integer programs of even larger size which are therefore typically more difficult to solve. The second objective is thereby provided with a weighting factor  $\gamma > 0$ . For  $\gamma \leq 1/n$  it was shown that a solution of the TBPP-FU always solves the TBPP as well and thus both problems are relatively close, see [Aydin, Muter, & Ilker Birbil \(2020\)](#). For other choices of the weighting parameter

- a solution with minimum number of servers (in terms of TBPP) does not have to be optimal for the TBPP-FU, see [Aydin et al. \(2020, Example 2.2\)](#),
- the possibility to decompose an instance (of the TBPP-FU) in a temporal sense is lost, see [Martinovic & Strasdat \(2022, Theorem 3\)](#),

so that a straightforward relation between the two problems does not exist anymore. Although the compact models for the TBPP-FU (called M1 and M2), originally introduced in [Aydin et al. \(2020\)](#), have been substantially improved over the past two years ([Martinovic & Strasdat, 2022; Martinovic, Strasdat, & Selch, 2021; Martinovic, Strasdat, Valério de Carvalho, & Furini, 2022](#)), only about 66% of the problem-specific benchmark instances can be solved optimally in reasonable time. In addition, if the benchmark sets formerly designed in [Dell'Amico et al. \(2020\)](#) for the traditional TBPP are now also taken into account, many more instances of moderate size cannot be tackled successfully, see [Martinovic, Strasdat, & Selch \(2021\)](#).

Thus, for both variants of temporal bin packing considered here, good exact approaches (either compact models or branch-and-price) have been found and their numerical performance has been optimized to a large extent, but numerous instances still remain unsolved. This article therefore proposes the concept of flow formulations, which is still (almost) entirely unexplored for both problems under consideration. Flow formulations form a powerful tool in cutting and packing, as they combine important structural properties (e.g., a good LP relaxation) with a large illustrativeness and a generally manageable model size, and so they can be handled efficiently by commercial solvers. In particular, extensions of the flow models originally described in detail for the first time in [Valério de Carvalho \(2002\)](#) have therefore been widely used in the recent past to present competitive exact approaches to fundamental optimization problems such as the cutting stock problem ([de Lima, Iori, & Miyazawa, 2022b; Delorme & Iori, 2020](#)), the skiving stock problem ([Martinovic, Delorme, Iori, Scheithauer, & Strasdat, 2020](#)), or the multiple knapsack problem ([Dell'Amico, Delorme, Iori, & Martello, 2019](#)). In particular, the importance of the

general methodology is also highlighted by the recent survey article ([de Lima, Alves, Clautiaux, Iori, & Valério de Carvalho, 2022a](#)).

### 1.3. Our Contribution

While all these very successful approaches have in common that they require a (pseudo-)polynomial number of states (nodes) and transitions (arcs) and thus allow the efficient treatment as an ILP formulation, such a graph-theoretic formulation for the TBPP is not yet known and not within reach. As already described for the example of the TKP, see [Arkin & Silverberg \(1987\)](#) and [Caprara et al. \(2013\)](#), the only way out is therefore via a graph which has an exponential number of states and transitions. Such an approach is also called a *combinatorial flow model* and is, however, according to the previous sources (and also [Furini, 2011, page 22](#)), only useful if, for example, the number of simultaneously active jobs at any point in time is very restricted – a property that is generally not given for the benchmark instances mentioned before. In addition, even the authors of a very recent work on exponential-size networks to tackle the TKP, see [Clautiaux, Detienne, & Guillot \(2021\)](#), admit that a straightforward application of the graph-theoretic idea does not lead to an efficient solution framework. Probably for these reasons, such an approach to temporal bin packing problems has not yet been investigated in the literature at all.

With this paper, we would like to contribute to foster the research on flow-based approaches to the two optimization problems under consideration. The main results of this work are the following:

- For both, the TBPP and the TBPP-FU, we present a layer-based combinatorial flow model. Here, each layer corresponds to a maximal clique of the interval graph belonging to the instance.
- Our approach is different from the previous attempts from the literature, see [Caprara et al. \(2013\)](#) and [Clautiaux et al. \(2021\)](#), as it uses another interpretation of states and transitions, leading to much smaller (but still exponentially large) networks.
- We improve the combinatorial flow models obtained in this way by valid inequalities.
- For both problems, the TBPP and the TBPP-FU, all known benchmark instances can be solved exactly in reasonable time. Moreover, significantly better computation times are achieved for those instances that could already be handled with the methods from the literature. As an outlook, we also try to explore the limits of our combinatorial arcflow model by dealing with instance sizes much larger than reported in [Dell'Amico et al. \(2020\)](#) and [Martinovic, Strasdat, & Selch \(2021\)](#).

We highlight that, although the paper just addresses two important application problems, which are related but have relevant structural differences, combinatorial flow models have a much wider applicability, and pave the way for further very powerful solution techniques to other interval scheduling problems discussed in the introductory parts.

## 2. The Temporal Bin Packing Problem: Preliminaries and Solution Methods

Let us consider a list of  $n \in \mathbb{N}$  items (jobs), specified by an item size (resource demand)  $c_i > 0$  and an activity interval (lifespan)  $[s_i, e_i]$  with  $s_i < e_i$ ,  $i \in I := \{1, \dots, n\}$ , and a sufficiently large number of homogeneous bins (servers) of capacity  $C > 0$ . We will refer to  $s_i$  and  $e_i$  by the *starting time* and *ending time* (or *terminating time*), respectively. Without loss of generality, we make the following assumptions:

- All input data are integers.
- The items are sorted with respect to non-decreasing starting times (where ties are broken in an arbitrary way).

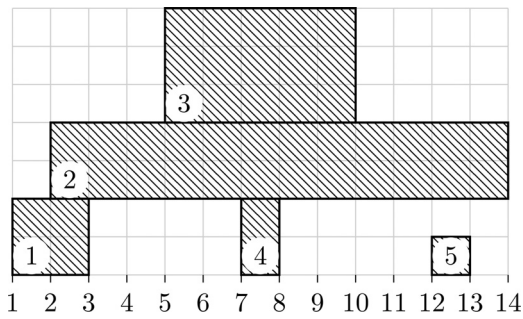


Fig. 1. An illustration of the instance  $E_0$ . The horizontal axis specifies the time instants, while the vertical grid measures the item sizes.

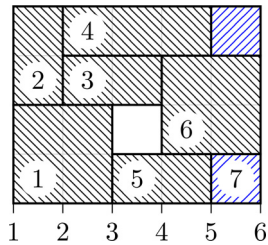


Fig. 2. An assignment of seven items to one bin of size  $C = 5$ , following an idea from bandwidth allocation in wireless networks, see Chen et al. (2002). The blue item  $i = 7$  with  $[s_i, e_i] = [5, 6]$  and  $c_i = 2$  is not placed in a connected manner. Moreover, there is no rearrangement of the items that enables representing item  $i = 7$  as a single rectangle of size  $1 \times 2$ , without destroying the rectangular structure of another item. In particular, this packing would not be feasible for the strip packing problem (SPP). (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

- The statement  $c_i \leq C$  holds for all  $i \in I$  (because the problem would become infeasible otherwise).

Then, the *temporal bin packing problem (TBPP)* requires to schedule the jobs to a minimum number of servers, so that the capacity of any server is respected at any instant of time. To briefly refer to a particular TBPP, we introduce the following well-known term.

**Definition 1.** A tuple  $E = (n, C, \mathbf{c}, \mathbf{s}, \mathbf{e})$ , where  $\mathbf{c}$ ,  $\mathbf{s}$ , and  $\mathbf{e}$  are  $n$ -dimensional vectors collecting the input-data (size, starting time, ending time) of the items, is called an *instance* (of the TBPP).

Typically, we refer to the set of time instants by  $T := \bigcup_{i \in I} \{s_i, e_i\}$ , and address the set of starting times by  $T_S := \bigcup_{i \in I} \{s_i\}$ . Moreover, the set  $I_t := \{i \in I \mid t \in [s_i, e_i]\}$  collects all jobs that are active at time  $t \in T$ .

**Example 1.** Let us consider the instance

$$E_0 = (5, 5, (2, 2, 3, 2, 1), (1, 2, 5, 7, 12), (3, 14, 10, 8, 13))$$

which is displayed in Fig. 1 and taken from Furini (2011, Section 3.2). Here, we have  $T = \{1, 2, 3, 5, 7, 8, 10, 12, 13, 14\}$ ,  $T_S = \{1, 2, 5, 7, 12\}$ , and (by way of example)  $I_7 = \{2, 3, 4\}$ . Note that an optimal solution requires two servers.

Notice that, although the items are visualized as ordinary rectangles in the “capacity-time plane” in Fig. 1, in a feasible solution of the TBPP they can in fact also be packed in such a way that they do not represent connected objects, see Fig. 2 for an example.

**Remark 1.** A similar illustration was already presented in Dell’Amico et al. (2020, Fig. 2) to show that the TBPP uses a different concept of feasible configurations than, for instance, the SPP. It should be noted that the instance depicted in Fig. 2 is smaller than the existing counterexample from the literature, both in terms of the number of items and the server capacity.

Following the structure of Kantorovich-type models for the BPP, see Kantorovich (1939), a first compact formulation for the TBPP was proposed in Dell’Amico et al. (2020) and can be obtained as follows. With  $K := \{1, \dots, n\}$  denoting the set of all servers, we can introduce two types of binary variables:

- We define  $z_k \in \{0, 1\}$  with  $z_k = 1$  if and only if server  $k \in K$  is used.
- We define  $x_{ik} \in \{0, 1\}$  with  $x_{ik} = 1$  if and only if item  $i \in I$  is assigned to server  $k \in K$ .

Then, we obtain the

**Compact Model for the TBPP (from Dell’Amico et al., 2020)**

$$z^{com} = \sum_{k \in K} z_k \rightarrow \min$$

$$\text{s.t. } \sum_{k \in K} x_{ik} = 1, \quad i \in I, \tag{1}$$

$$\sum_{i \in I_t} c_i x_{ik} \leq C \cdot z_k, \quad t \in T, k \in K, \tag{2}$$

$$x_{ik} \in \{0, 1\}, \quad i \in I, k \in K, \tag{3}$$

$$z_k \in \{0, 1\}, \quad k \in K. \tag{4}$$

The objective function minimizes the total number of servers in use. Moreover, the two sets of constraints make sure that any job is executed precisely once (see (1)) and that the capacity of the servers is respected at any instant of time (see (2)). Additionally, the latter prevent jobs from being assigned to unused servers at all.

For a fixed server  $k \in K$ , it is sufficient to require Conditions (2) only for all  $t \in T_S$ , since the load on a server can increase only at precisely these points in time. In fact, we can even go one step further.

**Definition 2 (Dell’Amico et al., 2020).** Let  $E$  be an instance of the TBPP and let  $t_1 < t_2 \in T_S$  follow each other directly in the chronologically ordered set  $T$  of all time instants. If  $t_2$  is not also an end time, then  $t_1$  is *dominated by*  $t_2$ . The set of all non-dominated starting times is referred to as  $T_S^{nd} \subseteq T_S$ .

In a situation like the one described in the definition, all jobs that are active at  $t_1$  are still active at  $t_2$ , meaning that the associated capacity condition (2) for  $t = t_2$  contains all the terms that would appear in the constraint for  $t = t_1$ , so that it dominates that restriction for all  $k \in K$ . Hence, Constraints (2) only need to be formulated for the non-dominated starting times.

**Example 2.** For our toy instance  $E_0$ , illustrated in Fig. 1, we conclude that  $T_S^{nd} = \{2, 7, 12\}$ . By way of example, the former element  $t_1 = 1 \in T_S$  is dominated by  $t_2 = 2$  and  $t_4 = 5$  by  $t_5 = 7$ .

Despite these possible improvements, such assignment-based ILP models usually have two major drawbacks, which are also evident here:

- The set of feasible solutions is highly symmetric due to permutations of the server indices.
- The LP bound coincides with a problem-specific generalization of the material bound, that is,  $\lceil \sigma / C \rceil$  where  $\sigma = \max_{t \in T} \sum_{i \in I_t} c_i$ . This bound is known to be rather poor, see Dell’Amico et al. (2020, Property 1) for a general observation and Dell’Amico et al. (2020, Fig. 3) for an exemplary instance needing strictly more bins than proposed by the lower bound.

Due to these disadvantageous properties of the compact model, a pattern-based approach has been established in the literature, see Dell’Amico et al. (2020), with a structure strongly reminiscent of the Gilmore-Gomory model of one-dimensional cutting stock problems, see Gilmore & Gomory (1961).

**Definition 3.** Any feasible assignment of jobs to a single server is called a *pattern*.

Mathematically, a pattern can be described by an  $n$ -dimensional incidence vector  $\mathbf{a} \in \{0, 1\}^n$  (or, equivalently, as a subset  $U \subseteq I$ ) the components of which contain the information whether item  $i \in I$  is part of the pattern ( $a_i = 1$ ) or not ( $a_i = 0$ ). Hence, the set of patterns for the TBPP is given by

$$\mathcal{P} := \mathcal{P}(E) := \left\{ \mathbf{a} \in \{0, 1\}^n \mid \sum_{i \in I_t} c_i a_i \leq C, t \in T \right\},$$

where  $T$  can also be replaced by  $T_s^{nd}$ . Due to the numerous combination possibilities, the cardinality of this pattern set typically grows exponentially with the number  $n$  of jobs appearing in an instance. Let  $\mathcal{J}$  denote an index set of  $\mathcal{P}$ , then we can introduce a decision variable  $\xi_j \in \{0, 1\}$  for each pattern  $j \in \mathcal{J}$ , stating whether it is used ( $\xi_j = 1$ ) or not ( $\xi_j = 0$ ). By that, we obtain the

**Exponential-size Model for the TBPP (from Dell’Amico et al., 2020)**

$$z^{exp} = \sum_{j \in \mathcal{J}} \xi_j \rightarrow \min$$

$$\text{s.t. } \sum_{j \in \mathcal{J}} \xi_j a_i^j = 1, \quad i \in I, \tag{5}$$

$$\xi_j \in \{0, 1\}, \quad j \in \mathcal{J}. \tag{6}$$

Again, the total number of servers is minimized while ensuring that any job is contained in precisely one pattern used, see Constraints (5). The exponential-size formulation does not contain any symmetry, and also its LP bound is generally better than that of the compact model presented before, see Dell’Amico et al. (2020, Property 4). Since the LP relaxation of the exponential-size model can be solved efficiently by column generation, the currently best solution approach for the TBPP, called B&P<sup>+</sup>, uses a branch-and-price algorithm based on that formulation, see Dell’Amico et al. (2020, Section 6). Before starting the actual (and costly) branch-and-price main procedure, the algorithm first tries to solve a given instance exactly using various lower bounds and heuristics.

- Determining appropriate lower bounds is mainly done by computing the rounded-up LP values  $\lceil z_{LP}^{exp,*} \rceil$ .
- During the first phases of the algorithm, the previously determined lower bounds are compared with a plethora of heuristic values. The heuristics used for this comparison are sorted by ascending difficulty and complexity. Thus, first an attempt is made to prove optimality for a given instance using very simple heuristics (e.g., first-fit techniques for the original and the lifted instance, see also Boschetti, Hadjiconstantinou, & Mingozi (2002), Clautiaux, Carlier, & Moukrim (2007) for the general concept of *lifting*), before moving to successively more sophisticated procedures culminating in token-based *diving heuristics*. The latter were proposed and discussed intensively in Sadykov, Vanderbeck, Pessoa, Tahiri, & Uchoa (2019) and intend to descend within a small part of the branch-and-bound tree according to a fixed (simple) heuristic rule until a suitable integer feasible solution is obtained. The value of the token thereby regulates that the numerically more difficult branching path  $\xi_j = 0$  (i.e., the decision not to use a certain pattern) can be chosen

only very rarely. In contrast, the simpler branching path  $\xi_j = 1$  (which allows continuing with a reduced and thus easier TBPP by removing the items occurring in the chosen pattern) may be used as often as desired.

Only when these previous techniques could not yet find a proven optimal solution for the given instance, the actual branch-and-price procedure is started. Here, the authors deviate from the classical pattern-based branching scheme, i.e.  $\xi_j = 0$  vs.  $\xi_j = 1$ . Instead, the algorithm uses a branching rule due to Ryan and Foster, see Ryan & Foster (1981), a strategy that considers several variables at once per node and, thus, typically offers rather balanced branching trees and a more efficient performance of the overall algorithm. For corresponding justifications and further explanations of the incorporation of the additional conditions into the respective subproblems, we refer to the general explanations in Barnhart, Johnson, Nemhauser, Savelsbergh, & Vance (1998, Section 4) as well as the problem-specific contributions to the TBPP in Dell’Amico et al. (2020, Section 6).

Altogether, the overall state-of-the-art algorithm B&P<sup>+</sup> is able to solve to proven optimality TBPP instances with up to 500 items and 150 non-dominated starting times, in reasonable computing times, as literally reported in Dell’Amico et al. (2020). However, not all the benchmark instances can be solved to proven optimality yet.

**3. A Combinatorial Flow-based Formulation for the Temporal Bin Packing Problem**

Although there is no graph-theoretical formulation for the TBPP in the relevant literature so far, two main concepts for the underlying TKP could constitute a starting point for further considerations. These approaches each describe layer-based graphs of exponential size, but they differ significantly in how such a layer is constructed. More specifically, the details of these two frameworks are given as follows:

- (1) In Caprara et al. (2013, Section 3.3), a *clique-based* idea already partly outlined in Arkin & Silverberg (1987) is discussed for the TKP. In that approach, the number of layers in the graph is determined by the number of maximal cliques of the interval graph belonging to the given instance. We note that there is a one-to-one relationship between the non-dominated starting times, see Definition 2, and the maximal cliques, implying a natural order among the maximal cliques. Moreover, the latter can be efficiently determined in polynomial time, see Biedl (2005) or Furini (2011, Algorithm 1) for an implementation with  $\mathcal{O}(n^2)$  time. Let us define  $c_0 := \emptyset$  and  $\nu_0 := \{\emptyset\}$  to represent an artificial first layer. For any of the remaining layers, consider a fixed maximal clique  $C_l$ . Then, the basic idea used in Caprara et al. (2013) is to define the states (nodes)  $\nu_l$  occurring in layer  $l$  of the graph as all feasible server allocations that can be built with the items of  $C_l$ , i.e., we have

$$\nu_l := \left\{ J \subseteq C_l : \sum_{i \in J} c_i \leq C \right\}.$$

In other words,  $\nu_l$  somewhat collects the “subpatterns” relevant for clique  $C_l$ . For any layer index  $l \geq 1$ , the arc set  $\mathcal{E}_l$  between layer  $l - 1$  and layer  $l$  is defined as follows:

$$(\underline{J}, \bar{J}) \in \mathcal{E}_l \iff \underline{J} \in \nu_{l-1}, \bar{J} \in \nu_l, \underline{J} \cap C_l = \bar{J} \cap C_{l-1}.$$

This definition particularly implies that both states,  $\underline{J}$  and  $\bar{J}$ , have to contain the same items from the set  $C_{l-1} \cap C_l$ , so that, among others, the artificial source node in  $c_0$  is connected to any node of the first layer. It is straightforward to see that there are at most  $\mathcal{O}(2^{|C_l|})$  nodes in layer  $l$ , so that the overall graph has an exponential size.

(II) In Clautiaux et al. (2021), on the other hand, an *event-based* approach is presented that leads to a graph whose number of layers is, by and large, determined by twice the number of items. Each of these layers belongs to a particular event, i.e., either the start or the completion of a job. In addition, there is an artificial last layer (with index  $2n + 1$ ), which consists only of a dummy sink node. In the graph itself, a state is described by a triple  $(e, c, \mathbf{a})$ , where  $e$  specifies the event and  $c$  is the capacity consumed by pattern  $\mathbf{a} \in \mathcal{P}$ . Of course, it is important to note that  $\mathbf{a}$  can only use those items that are consistent with the event under consideration. A transition from one state to the other can occur exactly when the associated item is picked up (at its start time) or released (at its end time). In the first case, it must also be checked whether the capacity condition is still fulfilled, while the second case requires the completed item to be part of the server state before. It is clear that, also for this variant, the state space will grow exponentially with the number of items.

Both approaches generally lead to very large graphs, which result, among other things, from focusing on the pattern set and/or having to represent each item twice (that is, by two events) during graph generation. However, the latter, as well as a certain degree of redundancy in the representation of nodes, is necessary for the application of the relaxation techniques presented in Clautiaux et al. (2021).

Therefore, in the following we would like to describe a possibility which, on the one hand, avoids redundant information in the labels and, on the other hand, leads to significantly smaller graphs, although we also put the patterns themselves in the center of our construction. In turn, however, our graph will sometimes have multiple arcs between two nodes. Let us start with two motivating examples to draw a direct connection to the previous attempts from the relevant literature, before going through the precise construction.

**Example 3.** We consider again the instance  $E_0$  from Example 1 and follow the idea of Caprara et al. (2013). First, the instance has three maximal cliques, namely  $C_1 = \{1, 2\}$ ,  $C_2 = \{2, 3, 4\}$ , and  $C_3 = \{2, 5\}$ , so the construction presented in Caprara et al. (2013) will lead to a graph consisting of four layers (three for the cliques and one dummy layer for the source node). For example, the first actual layer  $l = 1$  contains the states  $\mathcal{V}_1 = \{\emptyset, \{1\}, \{2\}, \{1, 2\}\}$ , while  $\mathcal{V}_2$  and  $\mathcal{V}_3$  consist of seven and four elements, respectively. Thus, in total, the graph has 16 nodes, which are connected by 32 arcs. Given its relatively large size we omit an illustration here, but note that one can be found in Furini (2011, Fig. 3.3). Looking at this graph a bit more closely, we see, for example, that the nodes  $\emptyset$  and  $\{1\}$  from layer  $l = 1$  are connected to exactly the same four nodes from layer  $l = 2$  (namely,  $\emptyset, \{3\}, \{4\}$ , and  $\{3, 4\}$ ). Thus, from the point of view of the arcs emanating from layer  $l = 1$ , these two states are to be evaluated as equivalent. This is also not unexpected, because at the end of the time interval relevant for clique  $C_1$ , job  $i = 1$  has already ended, so that it is irrelevant for the further server utilization (in the subsequent layers) which of the states  $\emptyset$  or  $\{1\}$  was once selected in layer  $l = 1$ . Hence, we suggest equipping a clique-based layer only with one representative per set of equivalent states (in the sense described before). Of course, then we have to attach the information which items were actually chosen to some other component of the graph. For that purpose, we will finally add a corresponding label to the arcs and allow multiple arcs between the states of two consecutive layers, if required. In Fig. 3, the combinatorial arcflow graph for  $E_0$  is depicted. Although the precise construction details have not yet been revealed, we see that it contains the same four clique-based layers, but only a total number of six nodes (illustrated as rectangles with rounded corners) and 15 arcs (black lines with rectangular label placed in

their center). So, in fact, the arcs of our approach somewhat carry the information of the states appearing in Furini (2011, Fig. 3.3), so that both, the number of nodes and arcs in Fig. 3, is much smaller than before.

**Example 4.** We highlight that, for the toy instance  $E_1 := (3, 8, (4, 4, 8), (0, 1, 3), (2, 4, 5))$  appearing in Clautiaux et al. (2021), our idea would lead to a graph having three layers (the two maximal cliques are  $C_1 = \{1, 2\}$  and  $C_2 = \{2, 3\}$ ), four nodes, and seven arcs, see Fig. 4. In contrast, the event-based approach also results in a much larger network consisting of seven layers, 15 nodes, and 18 arcs, see Clautiaux et al. (2021, Fig. 2).

Already from these examples we can see that our approach will lead to a much more efficient representation of the graphs. Although, from a theoretical point of view, they still are exponential in size, we will observe later in the numerical test calculations (see Sections 5 and 6) that both the generation of these graphs and the direct application of an ILP solver to the flow problem belonging to it does typically not consume an unreasonable amount of time anymore. In other words, for the first time our construction makes the exponential size controllable without requiring additional techniques.

We would now like to formalize the idea outlined in the previous examples. To this end, let  $E$  be a fixed instance of the TBPP, and let  $\mathcal{C} := \{C_1, \dots, C_m\}$  denote the set of its maximal cliques sorted in increasing order with respect to the associated non-dominated starting times. Moreover, we define two index sets  $L := \{1, \dots, m\}$  and  $L_0 := \{0\} \cup L$  to refer to the cliques and the layers, respectively. Lastly, remember that  $\mathcal{P}$  represents the set of all feasible patterns of  $E$ . In a slight abuse of notation, when describing and visualizing the ideas of our graph construction we will usually not refer to patterns by their incidence vectors. Instead, to not display too many redundant zero entries, we will make use of the corresponding subsets  $J \subseteq I$ . Since there is a one-to-one relation between these two concepts, no harm will arise from statements like  $J \in \mathcal{P}$ .

In our construction of the directed graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , any node will be referred to as a pair  $(l, J)$  with  $l \in L_0$  and some  $J \subseteq I$ . So, the first entry contains information about the layer, whereas the second specifies a subset of the items. Similarly, we will define an arc by a tuple  $(l - 1, l, J)$  with  $l \in L$  and some  $J \subseteq I$ , thus describing a transition from layer  $l - 1$  to  $l$  caused by  $J$ . Note that, whenever an arc is concerned, the subset  $J$  appearing in the tuple  $(l - 1, l, J)$  can also be interpreted as a (sub)pattern from

$$\mathcal{P}_l := \left\{ U \subseteq C_l : \sum_{i \in U} c_i \leq C \right\}.$$

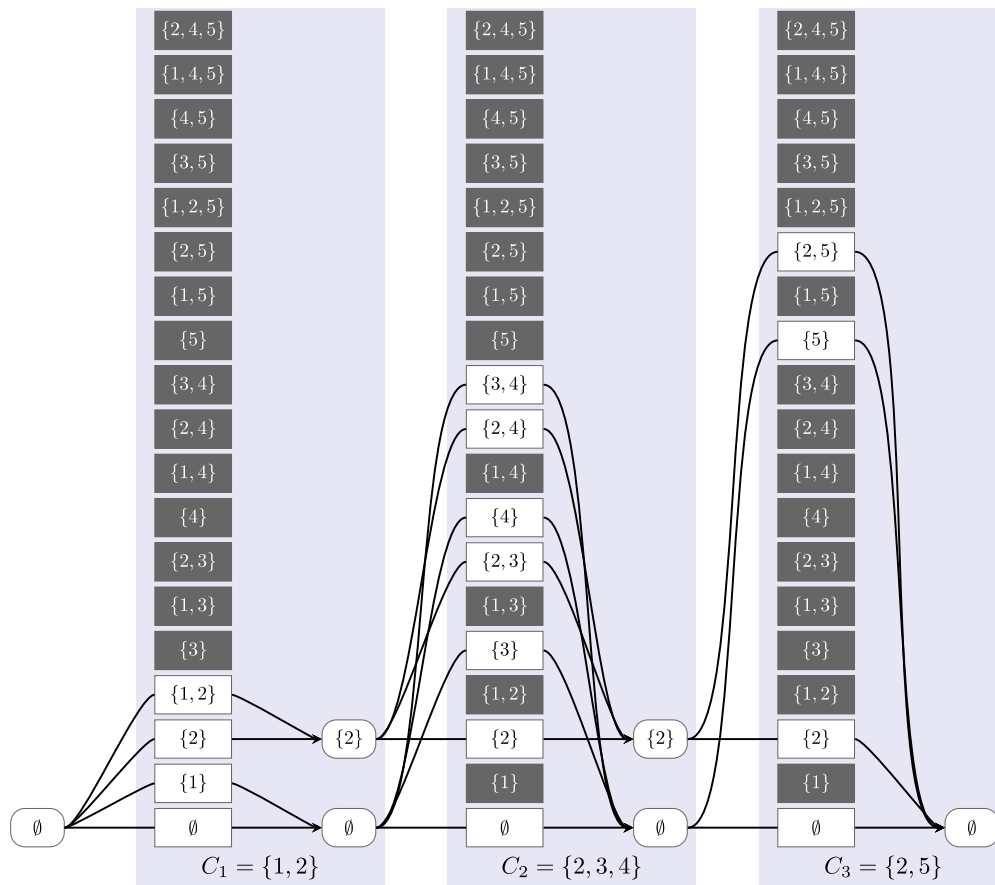
To initialize the set of nodes, let us define the dummy layer  $C_0$  with node set  $\mathcal{V}_0 := \{(0, \emptyset)\}$ . For any  $l \geq 1$ , we define (in two equivalent ways)

$$\mathcal{V}_l := \{(l, U \cap C_l \cap C_{l+1}) : U \in \mathcal{P}\} = \{(l, U \cap C_{l+1}) : U \in \mathcal{P}_l\}.$$

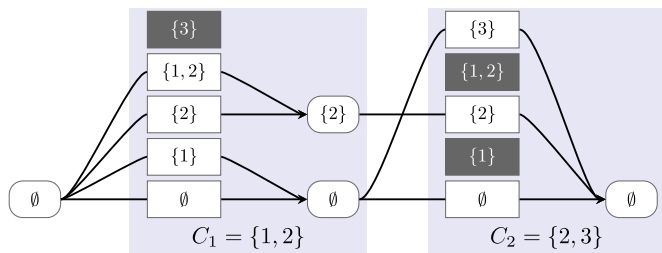
For the special case  $l = m$  (that is, the final layer), we use  $C_{m+1} := \emptyset$  in the above definition. So, in fact,  $\mathcal{V}_l$  collects the possible server states at time

$$e(l) := \max \{e_i : i \in C_l \setminus C_{l+1}\},$$

that can be observed if only the items exclusive to the set of cliques  $C_1, \dots, C_l$  processed so far are considered. In other words, these are the representative states that have been described in Example 3. Note that, due to the second possibility to define  $\mathcal{V}_l$ , we typically do not have to cope with the complete pattern set, but only with  $\mathcal{P}_l$ , so that the total number of vertices does not grow as fast as in the approaches known from the literature.



**Fig. 3.** The combinatorial arcflow graph for  $E_0$  consisting of four layers, six nodes, and 15 arcs. In contrast to the idea presented in Caprara et al. (2013), the arc and node labels carry some information relevant to the respective clique (so that, in our visualization, both of them are related to the blue background indicating a specific maximal clique). For the sake of completeness and to better understand the overall construction process, we depicted the entire pattern set  $\mathcal{P}$  in any layer, but cancelled out the elements that are infeasible because they do not use the items allowed for the respective clique. Those infeasible subpatterns are colored black and do not belong to any arc. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)



**Fig. 4.** The combinatorial arcflow graph for the instance  $E_1$  taken from Clautiaux et al. (2021).

Let  $\mathcal{E}_l$ ,  $l \in L$ , denote the set of arcs from layer  $l - 1$  to layer  $l$ . Similar (but not identical) to the approach presented in Caprara et al. (2013) and Furini (2011), there is some arc between two nodes  $(l - 1, J) \in \mathcal{V}_{l-1}$  and  $(l, \bar{J}) \in \mathcal{V}_l$  if and only if  $J \cap C_{l+1} = \bar{J} \cap C_{l-1}$ . However, for our network, we have to be more precise, because we can have multiple arcs between the same pair of nodes. So, the previous definition, in fact, just tells us that there is at least one arc between the two specific nodes, but further information is still missing. Hence, we demand that the nodes  $(l - 1, \underline{J}) \in \mathcal{V}_{l-1}$  and  $(l, \bar{J}) \in \mathcal{V}_l$  are connected if and only if

$$\exists J \in \mathcal{P}_l : J \cap C_{l-1} = \underline{J}, J \cap C_{l+1} = \bar{J}$$

holds, and we have to draw a separate arc for any such  $J$ . Hence, associated to each arc in  $\mathcal{E}_l$ , there is a unique subset of  $I$  (in fact, a

subpattern from  $\mathcal{P}_l$ ), called  $J$  in the previous definition, so that we can use this set to label the arc. As a consequence, we are allowed to abstractly refer to an arc by a tuple  $(l - 1, l, J)$  with some  $J \in \mathcal{P}_l$ , so that the notation introduced earlier is justified.

**Remark 2.** Note that any pattern  $J \in \mathcal{P}$  corresponds to a unique directed path from the source node  $(0, \emptyset)$  to the sink node  $(m, \emptyset)$ . Indeed, that path has to use the vertices  $(l, J \cap C_l \cap C_{l+1})$ ,  $l \in L_0$ , and the arcs  $(l - 1, l, J \cap C_l)$ ,  $l \in L$ . On the other hand, each path connecting the source and the sink node of the network via the arcs  $(0, 1, J_1)$ ,  $(1, 2, J_2)$ ,  $\dots$ ,  $(m - 1, m, J_m)$  defines a unique pattern, namely  $J = J_1 \cup \dots \cup J_m$  (or, more accurately, the corresponding incidence vector).

To conveniently formulate an integer optimization problem, let us collect all arcs referring to the bin allocation of item  $i \in I$  in the set  $\mathcal{E}(i)$ . In other terms, we define

$$\mathcal{E}(i) := \{(l - 1, l, J) \in \mathcal{E} : i \in J \setminus C_{l-1}, l \in L\}.$$

Moreover, the arcs entering and leaving a given state  $(l, J) \in \mathcal{V}$  will be denoted by  $\mathcal{E}^{in}(l, J)$  and  $\mathcal{E}^{out}(l, J)$ , respectively. Now, let us introduce an integer variable  $\xi_{l-1, l, J} \in \mathbb{Z}_+$  representing the units of flow carried by an arc  $(l - 1, l, J) \in \mathcal{E}_l$ ,  $l \in L$ . For the sake of simplicity, we will always use  $e$  to abbreviate the elements contained in a specific set of arcs. Then, we obtain the

**Combinatorial Arcflow Model for the TBPP**

$$z^{comb} = \sum_{e \in \mathcal{E}^{out}(0, \emptyset)} \xi_e \rightarrow \min$$

$$\text{s.t. } \sum_{e \in \mathcal{E}^{\text{in}}(l, \tilde{J})} \xi_e = \sum_{e \in \mathcal{E}^{\text{out}}(l, \tilde{J})} \xi_e, \quad (l, \tilde{J}) \in \mathcal{V} \setminus \{(0, \emptyset), (m, \emptyset)\}, \quad (7)$$

$$\sum_{e \in \mathcal{E}(i)} \xi_e = 1, \quad i \in I, \quad (8)$$

$$\xi_e \in \mathbb{Z}_+, \quad e \in \mathcal{E}_l, l \in L. \quad (9)$$

The objective function minimizes the total flow traversing the network (that is, the number of required servers), while Constraints (7) ensure the flow conservation at every vertex (except for the source and the sink node). Moreover, Conditions (8) manage that every job is executed precisely once. Note that, it is sufficient to restrict the flow to zero or one for most of the arcs. However, there are some arcs (namely, the arcs  $(l - 1, l, \emptyset)$ ,  $l \in L$ , representing a transition from one empty state to the next by means of  $J = \emptyset$ ) which can be used by multiple patterns. So, for exactly those arcs it is necessary to have integer-valued flow variables. To simplify the presentation of the model, here we do not differentiate between these two possibilities.

**Remark 3.** It is important to note that our construction does not require the items to be different in size, since the labelling just depends on the item indices. In the special case that items of the same size, in addition, have exactly the same activity interval, it would even be possible to use this information to slightly reduce the state space in the graph. Then the graph has to be constructed in a way to capture the number of items in each feasible state (in addition to their type). The reduction (compared to the original approach) comes from the fact that it is not necessary to construct all combinations of identical items. However, since most of the instances considered later (in Sections 5 and 6) either do not contain such “identical” items at all or only a very small number of them (relative to the total number of items), we will not examine this modification further here.

Now that we have presented and thoroughly explained the technical details of our graph construction, we would like to conclude by briefly discussing another example. In contrast to  $E_0$  and  $E_1$ , this one was not taken from the literature, but is already designed in such a way that we can then continue to work with it in the following section (when fire-ups have to be taken into account, too).

**Example 5.** Let us consider the instance  $E_2 := (5, 5, (1, 2, 5, 2, 4), (1, 2, 5, 6, 9), (7, 10, 6, 9, 10))$ , see Fig. A.1 in the appendix for a graphical illustration. The maximal cliques are given by  $C_1 = \{1, 2, 3\}$ ,  $C_2 = \{1, 2, 4\}$ , and  $C_3 = \{2, 5\}$ , so that  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  will consist of four layers, eight nodes, and 16 arcs, see Fig. 5. For the sake of completeness, in the appendix we also provide the graphs resulting from the approaches of Clautiaux et al. (2021) and Caprara et al. (2013) in Figs. A.2 and A.3, respectively.

#### 4. An Extension to the Temporal Bin Packing Problem with Fire-ups

##### 4.1. Preliminaries and a Basic Solution Approach

Minimizing the number of fire-ups and servers required is a very new aspect in the context of energy-efficient job-to-server scheduling introduced in Aydin et al. (2020) as the *temporal bin packing problem with fire-ups (TBPP-FU)*. The idea behind this is that not only the infrastructure as such but also its operating mode contributes to the energy consumption of the overall system. Consequently, this optimization problem assumes that servers that are temporarily unused can be put into a sleep mode or can be completely deactivated to save energy. However, such a server must

then be switched on again later, at the cost of a so-called *fire-up*. Both objectives are usually addressed by a weighted sum, scaling the number of fire-ups by some parameter  $\gamma > 0$ , in the objective function. This approach is also justifiable from a practical point of view, since both criteria more or less describe an energy consumption, and we can therefore bundle these objectives in a joint objective function.

**Definition 4.** A tuple  $E = (n, C, c, s, e, \gamma)$  where (in addition to the already known objects)  $\gamma > 0$  represents a scaling parameter is called an *instance* (of the TBPP-FU).

So far, the TBPP-FU has only been studied with respect to compact formulations. Two of these (called M1 and M2) were already suggested in Aydin et al. (2020), the article introducing this new optimization problem. Due to some serious drawbacks of these formulations, they were subject of numerous improvements over the past two years, see Martinovic, Strasdat, & Selch, 2021; Martinovic et al. (2022). The currently best formulation in the literature, see Martinovic & Strasdat (2022), has evolved from the original model M1. However, to keep the explanations short here we just introduce the basic version of M1 from Aydin et al. (2020). In addition to the classic assignment-based variables (that is,  $z_k$  and  $x_{ik}$ ) already known from Section 2, M1 uses the following two sets of variables to better access the temporal aspects of the problem under consideration:

- The decision variables  $y_{tk} \in \{0, 1\}$  will be interpreted in the sense that  $y_{tk} = 1$  represents a positive load on server  $k$  at time  $t$ .
- We use binary variables  $w_{tk} \in \{0, 1\}$  to state whether server  $k \in K$  was activated at time  $t \in T_S$ .

Then, we obtain the

**Assignment Model 1 (M1, original version from Aydin et al., 2020)**

$$z^{(1)} = \gamma \cdot \sum_{k \in K} \sum_{t \in T_S} w_{tk} + \sum_{k \in K} z_k \rightarrow \min$$

$$\text{s.t. } y_{tk} \leq \sum_{i \in I} c_i \cdot x_{ik} \leq y_{tk} \cdot C, \quad k \in K, t \in T, \quad (10)$$

$$\sum_{k \in K} x_{ik} = 1, \quad i \in I, \quad (11)$$

$$x_{ik} \leq y_{s_i, k}, \quad i \in I, k \in K, \quad (12)$$

$$y_{tk} \leq z_k, \quad k \in K, t \in T, \quad (13)$$

$$y_{tk} - y_{t-1, k} \leq w_{tk}, \quad k \in K, t \in T_S, \quad (14)$$

$$x_{ik} \in \{0, 1\}, \quad i \in I, k \in K, \quad (15)$$

$$y_{tk} \in \{0, 1\}, \quad k \in K, t \in T, \quad (16)$$

$$w_{tk} \geq 0, \quad k \in K, t \in T_S, \quad (17)$$

$$z_k \in \{0, 1\}, \quad k \in K. \quad (18)$$

The objective function collects the number of fire-ups (first sum) and the number of servers (second sum) and has to be minimized.



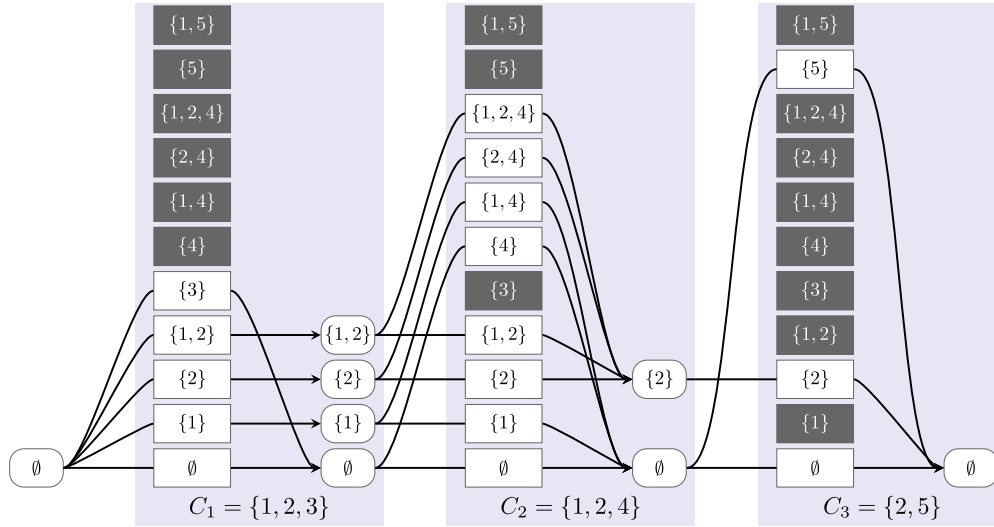


Fig. 5. The combinatorial arcflow graph for  $E_2$  from Example 5.

While Conditions (11) already appeared in the TBPP scenario, Constraints (10) manage that the capacity is respected whenever the considered server is active at the moment. In addition, Restrictions (12)–(14) are responsible for linking the different variable types consistently. Without going further into details, we mention that, meanwhile, this original version of M1 was improved by several general and problem-specific techniques like symmetry reduction, lifting, valid cuts, and heuristic-based information, see Martinovic & Strasdat (2022), Martinovic, Strasdat, & Selch, 2021, Martinovic et al. (2022).

#### 4.2. A Combinatorial Flow Formulation for the TBPP-FU

In the TBPP-FU a fire-up occurs if a server is used for the first time or if it is reactivated after it became inactive at some point back in time. So, basically, a fire-up should be registered whenever a server leaves the empty state. However, we cannot simply use the network introduced for the TBPP, because the new main challenge for our graph-theoretic approach is that no fire-up is necessary, if another job starts exactly at the time when a server intends to get empty. This means that it is necessary to differentiate between two possible empty states to count fire-ups correctly. More precisely, our construction will be based on a *true empty state* (shutdown), called  $\emptyset$ , and an *artificial empty state* (possible immediate resumption of activity), called  $\emptyset_A$ .

Before discussing the necessary changes in our graph construction, we have to identify in which situation the introduction of an additional empty state (in a specific layer) is mandatory. To this end, note that a server can only go into the artificial empty state if the last active job (on that server) runs exactly until another new job starts. Otherwise, the server is actually turned off and returns to the true empty state. To define the states accurately, we introduce the following notation for any clique  $C_l$ ,  $l \in L$ :

$$s(l) := \min \{s_i : i \in C_l \setminus C_{l-1}\},$$

$$e(l) := \max \{e_i : i \in C_l \setminus C_{l+1}\}.$$

Note that the second one already briefly appeared in Section 3, but we think mentioning it again will help to remember its meaning within the following constructions. From a descriptive point of view,  $s(l)$  is the earliest starting time of jobs that are introduced in  $C_l$ , whereas  $e(l)$  refers to the latest ending time of jobs that are completed in  $C_l$ . This means that a node representing the artificial empty state has to exist in layer  $l$ ,  $l \in L \setminus \{m\}$ , if and only if  $e(l) = s(l+1)$  holds.

**Example 6.** For the instance  $E_2$  from our previous example, we obtain the arcflow graph depicted in Fig. 6, if the TBPP-FU is considered.

Even if the essential differences to the graph generation of the TBPP have already been summarized, we would now like to discuss the formal definitions of the node and arc sets in more detail. To better distinguish between standard and artificial objects, we introduce additional tags in the description for both the states and the transitions. More precisely, we let  $\mathcal{V}_l^S$ ,  $l \in L_0$ , and  $\mathcal{V}^S$  denote the set of standard nodes, i.e., these are the vertices which already appeared in the TBPP graph from Section 3. In addition, we define

$$\mathcal{V}^A := \{(l, \emptyset_A) : l \in L \setminus \{m\}, e(l) = s(l+1)\}$$

to refer to the artificial empty states, whenever they are required in the respective layer of the graph. Of course, the set of all nodes is then given by  $\mathcal{V} := \mathcal{V}^S \cup \mathcal{V}^A$ .

To correctly reconstruct the fire-ups later, it is convenient to attach an additional fourth component to the description of an arc. First of all, we again have a set  $\mathcal{E}^S$  of standard arcs between two non-artificial states which can be inherited from the TBPP graph, see Section 3. However, any such arc is now referred to as a 4-tuple  $(l-1, l, J, S)$ , where the tag ‘S’ refers to ‘standard’. Similarly, any arc involving an artificial state can be assigned to one of the following three groups:

$$\mathcal{E}^{A\leftarrow} := \left\{ (l-1, l, J, \leftarrow) : l \in L, J \in \mathcal{P}_l, \max_{i \in J} \{e_i\} = s(l+1) \right\},$$

$$\mathcal{E}^{A\rightarrow} := \left\{ (l-1, l, J, \rightarrow) : l \in L, J \in \mathcal{P}_l, \min_{i \in J} \{s_i\} = e(l-1) \right\},$$

$$\mathcal{E}^{A\leftrightarrow} := \left\{ (l-1, l, J, \leftrightarrow) : (l-1, l, J, \rightarrow) \in \mathcal{E}^{A\rightarrow} \text{ and } (l-1, l, J, \leftarrow) \in \mathcal{E}^{A\leftarrow} \right\}.$$

We define the set of all arcs as  $\mathcal{E} := \mathcal{E}^S \cup \mathcal{E}^{A\leftarrow} \cup \mathcal{E}^{A\rightarrow} \cup \mathcal{E}^{A\leftrightarrow}$ . The new types of arcs represent the following transitions.

- An arc  $(l-1, l, J, \leftarrow) \in \mathcal{E}^{A\leftarrow}$  connects a state in  $\mathcal{V}_{l-1}^S$  (namely, the state belonging to  $J \cap C_{l-1}$ ) to the subsequent artificial empty state  $(l, \emptyset_A) \in \mathcal{V}^A$ . This corresponds to a server which would get empty at time  $s(l+1)$ , but could continue execution if a suitable job is assigned.
- Similarly, an arc  $(l-1, l, J, \rightarrow) \in \mathcal{E}^{A\rightarrow}$  starts in the artificial empty state  $(l-1, \emptyset_A) \in \mathcal{V}^A$  and goes to a state in  $\mathcal{V}_l^S$  (namely,

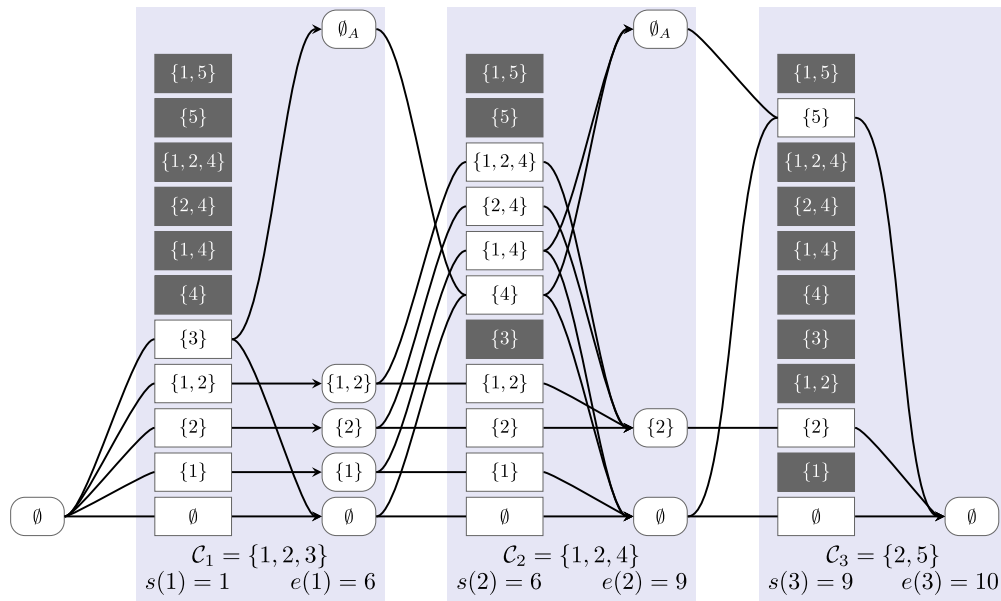


Fig. 6. The combinatorial arcflow graph for the TBPP-FU and the instance  $E_2$ .

the state belonging to  $J \cap C_l$ ). In this transition, a server which got empty recently is directly reactivated (without requiring a fire-up).

- An arc  $(l-1, l, J, \leftrightarrow) \in \mathcal{E}^{A \leftrightarrow}$  connects one artificial empty state  $(l-1, \emptyset_A) \in \mathcal{V}^A$  to the subsequent artificial empty state  $(l, \emptyset_A) \in \mathcal{V}^A$  (via the subpattern  $J \in \mathcal{P}_l$ ). In fact, this transition is a combination of the previous ones.

**Remark 4.** For the instance  $E_2$  dealt with in Section 3, we obtain the following sets of non-standard arcs

$$\begin{aligned} \mathcal{E}^{A \leftarrow} &= \{(0, 1, \{3\}, \leftarrow), (1, 2, \{1, 4\}, \leftarrow), (1, 2, \{4\}, \leftarrow)\}, \\ \mathcal{E}^{A \rightarrow} &= \{(1, 2, \{4\}, \rightarrow), (2, 3, \{5\}, \rightarrow)\}, \\ \mathcal{E}^{A \leftrightarrow} &= \{(1, 2, \{4\}, \leftrightarrow)\}, \end{aligned}$$

which can also be found in Fig. 6. Note that extending the representation of arcs to 4-tuples is indeed necessary, because in our example there are four different scenarios to move from layer  $l=1$  to layer  $l=2$  via the subpattern  $J = \{4\}$ , but the arcs  $(1, 2, \{4\}, \rightarrow)$ ,  $(1, 2, \{4\}, \leftrightarrow)$  do not contribute to a fire-up.

A potential issue with our definition of the states is that there is no longer a one-to-one relationship between patterns (i.e., the elements of  $\mathcal{P}$ ) and directed paths in the graph. As in the TBPP case, we still have that each path from the source to the sink corresponds to a feasible pattern, but there may be different paths leading to the same pattern, in general. More precisely, this is the case if the set  $\mathcal{V}^A$  is nonempty. In particular, for any path which goes through some artificial empty state  $(l, \emptyset_A) \in \mathcal{V}^A$ , there is another path using the node  $(l, \emptyset) \in \mathcal{V}^S$  instead and ends up with the same pattern, see Fig. A.4 in the appendix for an example. The necessary arcs for this replacement must exist by the definition of the graph. As regards our optimization model, however, this ambiguity is typically not problematic since a path through an artificial empty state is preferred over the equivalent path using a true empty state because of the fire-up costs  $\gamma > 0$ .

For the purpose of a preferably convenient modeling, we again let  $\mathcal{E}(i)$  denote the set of all arcs which represents the starting of job  $i \in L$ . Moreover, with  $\kappa$  acting as a generic tag symbol, we specify the incoming and outgoing arcs of some node  $(l, \tilde{J}) \in \mathcal{V}$  as follows:

low:

$$\begin{aligned} \mathcal{E}^{in}(l, \tilde{J}) &:= \begin{cases} \{(l-1, l, J, \kappa) \in \mathcal{E} : J \cap C_l = \tilde{J}, \\ \kappa \in \{S, \rightarrow\}\} & \text{if } (l, \tilde{J}) \in \mathcal{V}^S, \\ \{(l-1, l, J, \kappa) \in \mathcal{E} : \kappa \in \{\leftarrow, \leftrightarrow\}\} & \text{if } (l, \tilde{J}) = (l, \emptyset_A) \in \mathcal{V}^A. \end{cases} \\ \mathcal{E}^{out}(l, \tilde{J}) &:= \begin{cases} \{(l, l+1, J, \kappa) \in \mathcal{E} : J \cap C_{l+1} = \tilde{J}, \\ \kappa \in \{S, \leftarrow\}\} & \text{if } (l, \tilde{J}) \in \mathcal{V}^S, \\ \{(l, l+1, J, \kappa) \in \mathcal{E} : \kappa \in \{\rightarrow, \leftrightarrow\}\} & \text{if } (l, \tilde{J}) = (l, \emptyset_A) \in \mathcal{V}^A. \end{cases} \end{aligned}$$

The main difference to the previous combinatorial flow model is that we have to correctly perceive the fire-ups. To this end, we make use of set

$$\mathcal{E}^{FU} := \{(l-1, l, J, \kappa) \in \mathcal{E} : l \in L, J \neq \emptyset, J \cap C_{l-1} = \emptyset, \kappa \in \{S, \leftarrow\}\}.$$

In other words, a fire-up has to be counted when connecting the true empty state with some other state via a standard arc (with tag ‘S’) corresponding to some pattern  $J \neq \emptyset$  or, alternatively, with an artificial empty state via an arc from  $\mathcal{E}^{A \leftarrow}$ . In particular, the first possibility also takes care of recognizing the first activation of a server.

With all these ingredients at hand, we can introduce integer variables  $\xi_{l-1, l, J, \kappa} \in \mathbb{Z}_+$  denoting the units of flow carried by arc  $(l-1, l, J, \kappa) \in \mathcal{E}_l$ ,  $l \in L$ . Then, we obtain the

**Combinatorial Arcflow Model for the TBPP-FU**

$$\begin{aligned} z^{comb} &= \sum_{e \in \mathcal{E}^{out}(0, \emptyset)} \xi_e + \gamma \cdot \sum_{e \in \mathcal{E}^{FU}} \xi_e \rightarrow \min \\ \text{s.t. } \sum_{e \in \mathcal{E}^{in}(l, \tilde{J})} \xi_e &= \sum_{e \in \mathcal{E}^{out}(l, \tilde{J})} \xi_e, \quad (l, \tilde{J}) \in \mathcal{V} \setminus \{(0, \emptyset), (m, \emptyset)\}, \quad (19) \end{aligned}$$

$$\sum_{e \in \mathcal{E}(i)} \xi_e = 1, \quad i \in L, \quad (20)$$

$$\xi_e \in \mathbb{Z}_+, \quad e = (l-1, l, J, \kappa) \in \mathcal{E}_l, l \in L. \quad (21)$$

In fact, the general form and the interpretation of the constraints did not change when moving from the TBPP to the TBPP-FU. The only considerable difference is that the objective function now also contains a sum collecting the fire-up terms.

## 5. Numerical Tests: Preliminaries

### 5.1. Computational Environment and Test Instances

In the literature, the following benchmark sets (referred to as *Category (A)* and *Category (B)*) have been described for the two problems under consideration:

- (A) In [Aydin et al. \(2020, Section 5\)](#), the authors suggested 48 differently characterized groups of 5 instances each, forming a set of 240 instances in total. All instances share the values  $C = 100$ ,  $\gamma = 1$ , but any two groups differ in precisely one of the criteria:

- number of items:  $n \in \{50, 100, 150, 200, 500, 1000\}$ ,
- time horizon: dense scenario ( $\bar{s} := \max_{i \in I} \{s_i\} = n$ ) vs. relaxed scenario ( $\bar{s} = 1.2n$ ),
- job duration: short ( $'d_S'$ ) vs. long ( $'d_L'$ ),
- capacity consumption: low ( $'c_L'$ ) vs. high ( $'c_H'$ ).

For the precise construction, we refer the interested reader to the aforementioned publication. From the input data, we can see that there is a wide range of possible values, especially with respect to the number of jobs, and these can be used to further decompose (A) into two subcategories:

- (A1) In this subset, we would like to summarize those 160 instances with values  $n \leq 200$ . Those instances have been tackled and (partially) solved in [Aydin et al. \(2020\)](#) using exact approaches, and they are also used in the articles dealing with improved compact formulations, see [Martinovic, Strasdat, & Selch, 2021](#), [Martinovic et al. \(2022\)](#).
- (A2) In this subset, we would like to gather the 80 significantly more difficult instances with  $n \in \{500, 1000\}$ . These have been treated in the literature so far exclusively with heuristic methods, see [Aydin et al. \(2020\)](#), so no information about optimal solutions is available.
- (B) In [Dell'Amico et al. \(2020, Section 7\)](#) the authors introduced a set of 1500 instances with  $C = 100$ , originating from an earlier investigation of the TKP in [Caprara et al. \(2013\)](#). In contrast to *Category (A)*, the range of the item sizes is not that restricted, and the main input parameter is given by the number of non-dominated starting times (maximal cliques). More precisely, for any  $|T_S^{nd}| \in \{10, 20, 30, \dots, 150\}$  a set of 100 instances (divided into ten classes called I–X) is considered. Any class is described by a parameter  $a$  denoting the average number of items per clique, and a parameter  $b$  influencing the job duration. So, the higher the  $a$ -parameter or the lower the  $b$ -parameter of an instance, the more jobs will have to be assigned in total. The full details of that construction can be found in [Dell'Amico et al. \(2020\)](#). Here, we just highlight that Classes VI and IX possess relatively small  $b$ -parameters, whereas Classes VIII–X exhibit relatively large values of the  $a$ -parameter, so these classes contain the more challenging instances.

**Remark 5.** For the sake of completeness, it should be noted that *Category (B)* also contains two sets of 100 instances each with  $|T_S^{nd}| \in \{5, 15\}$ . However, according to the experiments conducted in [Dell'Amico et al. \(2020\)](#), we will not make use of these (relatively easy) instance sets.

In the following subsections, we compare our combinatorial arcflow approach with the best known solution methods from the literature. To recapitulate, we have the B&P<sup>+</sup> algorithm from [Dell'Amico et al. \(2020\)](#) for the TBPP and a compact model for the TBPP-FU (originating from the formulation M1) presented in [Martinovic & Strasdat \(2022\)](#). Following the relevant literature, the instances of *Category (B)* will be used for both problems under

consideration, whereas *Category (A)* is exclusive to the TBPP-FU. The new flow-based approach is coded in Python (version 3.10.1) and solved by Gurobi (version 9.5) on an AMD A10-5800K processor with 16 GB RAM. Unless stated otherwise, we use a time limit of 30 minutes. However, especially for some very large instances appearing in [Subsect. 6.3](#), we will also perform computations without any time limit.

### 5.2. Structural Comparison of the Graph-based Approaches

Before dealing with the concrete numerical performance of our new approach, we would like to study its general applicability in more detail. This is done in particular against the background of allowing a comparison to the approaches from the literature, but also to show that the exponential size of the network, which still exists (in theory), now appears to be much more controllable. To this end, we collected the average size of the three graphs (in terms of nodes and arcs) for two representative instance sets from *Category (B)* in [Table 1](#).

**Remark 6.** To give some more information about the instances, we also list the average size of the cliques  $|C_l|_{avg}$  (i.e., the items per clique) and the average total number of items  $|I|_{avg}$ . Note that, effectively, the first value is identical to the mean value of the  $a$ -parameter used to construct the instance class in *Category (B)*, see [Dell'Amico et al. \(2020\)](#) for the details. Later, in [Table 5](#), the same information will be shown for *Category (A)* as well.

We clearly see that combinatorial arcflow (termed 'CAF') is much smaller than the two competitors from the literature. Having a look at the average numbers of states and transitions, there are remarkable savings of between 85 and 90% compared to the event-based graph presented in [Clautiaux et al. \(2021\)](#). A similar, but not to the identical degree superior, result is obtained in comparison with the layer-based idea from [Caprara et al. \(2013\)](#). Here, the reduction in terms of nodes is slightly above 37%, whereas the number of arcs decreases by almost 60%.

**Remark 7.** An interesting side aspect is given by the observation that, as a direct consequence of the construction itself, the number of nodes in the approach from [Caprara et al. \(2013\)](#) is almost equal to the number of arcs in CAF. In fact, only the source node of [Caprara et al. \(2013\)](#) does not appear as an arc in our implementation.

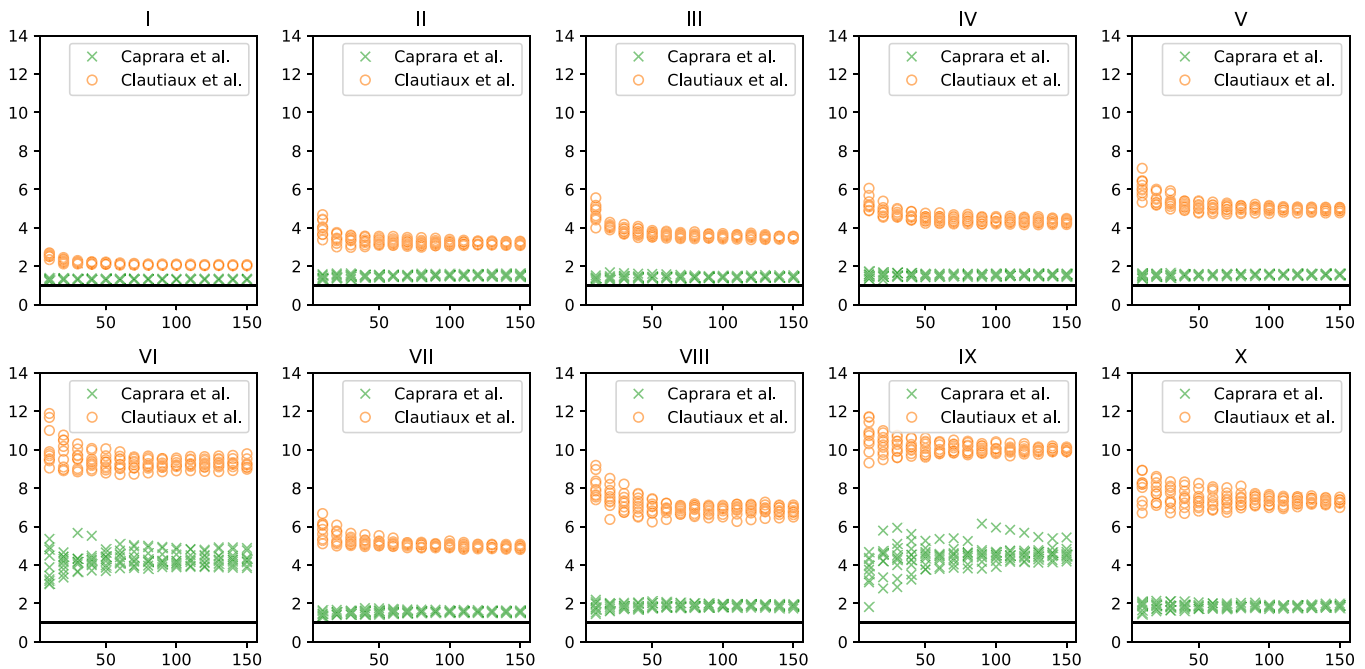
To also obtain a rough impression of the numerical data with respect to other parameter values of  $|T_S^{nd}|$  or the specific class index, in addition to the representative numbers in [Table 1](#), we display the normalized arc numbers (i.e., the value  $|\mathcal{E}|$  of CAF is scaled to 1) of *Category (B)* for the approaches from [Caprara et al. \(2013\)](#) and [Clautiaux et al. \(2021\)](#) in [Fig. 7](#). This is of particular interest because the number of arcs is identical to the number of integer variables in the ILP model and, thus, significantly determines the solution efforts. First, it is noticeable that after initial fluctuations, a relatively constant ratio is obtained for increasing values of  $|T_S^{nd}|$ . We attribute this to the fact that the construction of a fixed instance class, independent of  $|T_S^{nd}|$  itself, always follows the same principle, and, therefore, basic structural properties (like the "density" of the arcs between the layers) of the graphs are preserved for each of the three approaches. It can be further seen that even for values of  $|T_S^{nd}|$  other than those studied in [Table 1](#), CAF achieves a large saving over the networks from the literature, especially compared to [Clautiaux et al. \(2021\)](#). The effect varies by class and is particularly pronounced, then also in comparison with [Caprara et al. \(2013\)](#), for Classes VI and IX, which have the most items (see [Table 1](#)) and are, thus, particularly challenging.

However, all these observations are pointless if the generation of the still exponentially large network alone takes too much time

**Table 1**

The average number of nodes and arcs (in units of  $10^3$ ) for three different arcflow graphs: the event-based version from Clautiaux et al. (2021), the layer-based variant from Caprara et al. (2013), and combinatorial arcflow (“CAF”) presented in this work. For the sake of exposition, we just consider two representative sets of instances from Category (B).

$ T_5^{nd} $	Class	$ C_1 _{avg}$	$ I _{avg}$	number of nodes ( $ \mathcal{V} $ )			number of arcs ( $ \mathcal{E} $ )		
				Clautiaux et al. (2021)	Caprara et al. (2013)	CAF	Clautiaux et al. (2021)	Caprara et al. (2013)	CAF
50	I	10.0	59.0	4.7	2.5	<b>1.9</b>	5.2	3.2	<b>2.4</b>
	II	15.0	97.0	20.1	6.7	<b>4.9</b>	22.0	10.1	<b>6.7</b>
	III	20.0	110.0	33.3	9.8	<b>7.5</b>	35.6	14.0	<b>9.8</b>
	IV	25.0	139.1	113.0	26.9	<b>19.6</b>	120.4	41.6	<b>26.9</b>
	V	30.0	161.2	214.1	45.0	<b>32.6</b>	226.5	70.0	<b>45.0</b>
	VI	30.0	339.7	424.7	48.8	<b>22.2</b>	451.4	208.6	<b>48.8</b>
	VII	30.0	161.8	209.2	42.9	<b>31.4</b>	220.7	66.9	<b>42.9</b>
	VIII	30.0	213.5	330.0	49.6	<b>30.4</b>	349.3	93.9	<b>49.6</b>
	IX	29.9	354.4	442.4	46.8	<b>20.9</b>	468.3	208.7	<b>46.8</b>
	X	34.9	231.3	436.4	62.3	<b>40.8</b>	457.9	114.2	<b>62.3</b>
Average		25.5	186.7	222.8	34.1	<b>21.2</b>	235.7	83.1	<b>34.1</b>
100	I	10.0	109.0	8.5	4.6	<b>3.7</b>	9.5	6.1	<b>4.6</b>
	II	15.0	181.0	39.0	13.4	<b>9.8</b>	42.7	20.5	<b>13.4</b>
	III	20.0	201.3	61.6	18.6	<b>14.4</b>	65.9	26.6	<b>18.6</b>
	IV	25.0	255.1	209.4	51.0	<b>37.2</b>	223.2	79.1	<b>51.0</b>
	V	30.0	294.0	354.2	75.3	<b>54.8</b>	374.8	118.5	<b>75.3</b>
	VI	30.0	657.1	760.8	87.8	<b>40.7</b>	808.0	375.2	<b>87.8</b>
	VII	30.0	295.6	371.4	78.1	<b>57.4</b>	392.1	122.5	<b>78.1</b>
	VIII	30.0	402.5	630.3	96.7	<b>60.5</b>	666.6	181.8	<b>96.7</b>
	IX	30.1	689.9	872.8	92.1	<b>41.1</b>	923.9	422.3	<b>92.1</b>
	X	35.0	432.3	982.4	140.7	<b>91.2</b>	1032.0	258.0	<b>140.7</b>
Average		25.5	351.8	429.0	65.8	<b>41.1</b>	453.9	161.0	<b>65.8</b>



**Fig. 7.** Number of arcs (normalized to CAF) for the networks presented in Caprara et al. (2013) and Clautiaux et al. (2021) for the ten classes belonging to Category (B). The horizontal axis represents the parameter  $|T_5^{nd}|$  of the instances.

to be able to solve it subsequently. To reject this possible criticism as well, we list the average times to generate the graph (including the corresponding ILP model for the TBPP) in Table 2. Given this data and the time limit of 30 minutes, it can be stated that the construction of the graph generally requires only a relatively small amount of time, even for larger instances, and thus the ILP model can be easily passed to a solver. We also point out that with other programming languages, even shorter runtimes could be expected for model generation, because the performance of statically typed, compiled languages is typically much better (compared to Python).

**Remark 8.** It is interesting to note that it does not matter which of the networks (TBPP or TBPP-FU) is considered. The times to generate the graphs (and the corresponding ILP formulation) are very close to each other, because there is only a tiny difference with respect to the number of nodes and arcs between the two frameworks. Since this argument should be clear from the constructions described in the previous sections, in Tables 1 and 2 we do not intend to provide additional data for the case of the TBPP-FU.

As a conclusion of this subsection, we would like to summarize that our graph is, in fact, much smaller than the previous

**Table 2**  
The average time (in seconds) to construct the CAF network (and the associated ILP).

$ T_S^{nd} $	10	20	30	40	50	60	70	80	90	100	110	120	130	140	150
I	0.0	0.1	0.1	0.2	0.2	0.2	0.3	0.3	0.4	0.4	0.5	0.5	0.5	0.5	0.6
II	0.1	0.2	0.3	0.5	0.6	0.7	0.8	1.0	1.1	1.2	1.3	1.4	1.6	1.7	1.9
III	0.1	0.3	0.5	0.7	0.9	1.0	1.2	1.3	1.6	1.7	1.9	2.2	2.6	2.8	3.1
IV	0.5	1.0	1.3	2.0	2.5	3.0	3.4	3.8	4.3	4.8	5.4	5.8	6.4	6.6	7.0
V	0.7	1.4	2.3	3.5	4.3	4.8	5.4	6.0	6.5	7.1	8.2	9.1	9.9	10.7	11.4
VI	0.7	1.3	2.1	2.8	3.7	4.4	5.0	5.6	6.3	6.9	7.7	8.7	9.3	10.2	11.0
VII	1.0	1.8	2.6	3.3	4.1	4.7	5.5	6.2	6.9	7.5	8.5	9.2	10.0	10.9	11.9
VIII	0.7	1.5	2.2	3.3	4.3	5.3	6.5	7.2	8.0	8.7	9.8	10.3	11.1	12.2	12.8
IX	0.5	1.3	2.2	2.8	3.6	4.3	4.9	5.5	6.4	7.2	7.9	8.8	9.3	9.8	10.5
X	1.1	2.1	3.4	4.5	6.0	7.2	8.7	10.1	11.2	13.0	14.5	15.9	17.7	20.1	21.9
Average	0.6	1.1	1.7	2.3	3.0	3.6	4.2	4.7	5.3	5.8	6.6	7.2	7.8	8.6	9.2

**Table 3**  
Numerical comparison (for the TBPP) between CAF and the best solution approach from the literature, that is, B&P+ from Dell’Amico et al. (2020), for Category (B) (ordered by number of non-dominated starting times).

$ T_S^{nd} $	CAF with $t_{max} = 1800s$		Dell’Amico et al. (2020) with $t_{max} = 3600s$	
	t	opt	t	opt
10	<b>0.4</b>	<b>(100)</b>	2.0	(100)
20	<b>1.4</b>	<b>(100)</b>	5.2	(100)
30	<b>2.5</b>	<b>(100)</b>	9.3	(100)
40	<b>3.7</b>	<b>(100)</b>	15.7	(100)
50	<b>6.9</b>	<b>(100)</b>	27.9	(100)
60	<b>10.0</b>	<b>(100)</b>	63.9	(100)
70	<b>15.7</b>	<b>(100)</b>	115.1	(100)
80	<b>16.5</b>	<b>(100)</b>	132.6	(100)
90	<b>20.4</b>	<b>(100)</b>	151.8	(99)
100	<b>23.9</b>	<b>(100)</b>	168.8	(99)
110	<b>26.7</b>	<b>(100)</b>	218.0	(98)
120	<b>30.7</b>	<b>(100)</b>	237.9	(99)
130	<b>38.6</b>	<b>(100)</b>	432.9	(94)
140	<b>40.2</b>	<b>(100)</b>	475.1	(92)
150	<b>51.5</b>	<b>(100)</b>	559.4	(90)
Average	<b>19.3</b>	<b>(1500)</b>	174.4	(1471)

approaches from the literature, and its generation takes only a relatively short time (especially measured against the available time limit). The graph we present is therefore very well suited to be used for an exact solution by ILP solvers. The corresponding results obtained from numerical computations involving the benchmark sets presented before will be documented and discussed in the following subsections.

**6. Numerical Tests: Solving Benchmark Instances with CAF**

6.1. Numerical Results for the TBPP

As mentioned in Section 2, currently the best method for solving TBPP instances is the B&P+ algorithm developed in Dell’Amico et al. (2020). This approach has already been able to solve the vast majority of the 1500 TBPP-specific benchmark instances (that is, Category (B)) in reasonable time. However, there are still 29 of these instances – all of which having  $|T_S^{nd}| \geq 90$  and originating exclusively from the Classes V-X – where no proven optimal solution was found. To this end, in Tables 3 and 4 we copied the results of the previous state of the literature (so, the mean running times and number of optimal solutions reported in Dell’Amico et al., 2020), and compare them with the performance of CAF. We note that in Dell’Amico et al. (2020) a time limit of one hour was used.

We see that CAF can now solve every single instance, typically requiring significantly less computation time in each of the subsets of instances considered in the tables. Looking at the averages over

**Table 4**  
Numerical comparison (for the TBPP) between CAF and the best solution approach from the literature, that is, B&P+ from Dell’Amico et al. (2020), for Category (B) (ordered by instance classes).

Class	CAF with $t_{max} = 1800s$		Dell’Amico et al. (2020) with $t_{max} = 3600s$	
	t	opt	t	opt
I	<b>0.2</b>	<b>(150)</b>	0.6	(150)
II	<b>2.0</b>	<b>(150)</b>	9.4	(150)
III	<b>1.8</b>	<b>(150)</b>	8.1	(150)
IV	<b>25.3</b>	<b>(150)</b>	60.1	(150)
V	<b>49.0</b>	<b>(150)</b>	396.0	(146)
VI	<b>11.1</b>	<b>(150)</b>	214.7	(146)
VII	<b>48.9</b>	<b>(150)</b>	248.7	(147)
VIII	<b>17.7</b>	<b>(150)</b>	223.9	(144)
IX	<b>10.6</b>	<b>(150)</b>	82.3	(147)
X	<b>26.2</b>	<b>(150)</b>	500.3	(141)
Average	<b>19.3</b>	<b>(1500)</b>	174.4	(1471)

all instances (in the last row), we observe a reduction in computation time of almost 90%, which on closer inspection, for example, becomes even larger for particularly challenging choices of the parameter  $|T_S^{nd}|$ , see Table 3. On the other hand, we also note that for some of the more difficult instance classes, such as Classes VI, VIII, and X, there are still much more remarkable performance gains of up to roughly 95%, see Table 4. By way of example, for Class VI we already saw in Fig. 7 that, compared to the large number of items, CAF leads to a particularly small graph representation, so that the convincing performance noticed in Table 4 is consistent with the previous results.

**Remark 9.** We note that the tabulated computation times for CAF are only the pure solution times of the ILP solver. However, this does not distort the previous statements in any way, as we have seen in Table 2 how small the modeling times turn out to be despite the exponential size. So, even if we added these times to the average values of the solution time, we would still observe a clear victory of CAF for each considered subset.

A somewhat more detailed overview of both solution methods is shown in a performance profile, see Fig. 8. This illustration displays the percentage of optimally solved instances over time, and it reveals that, apart from a small interval around a computation time of one second, CAF is strictly better than B&P+ at any point in time.

6.2. Numerical Results for the TBPP-FU

Since we are now able to solve any known benchmark instance for the TBPP in a short time with the help of our new approach, we would now like to turn to the somewhat more complex TBPP-FU. At the beginning, we note that we have added the lower bound

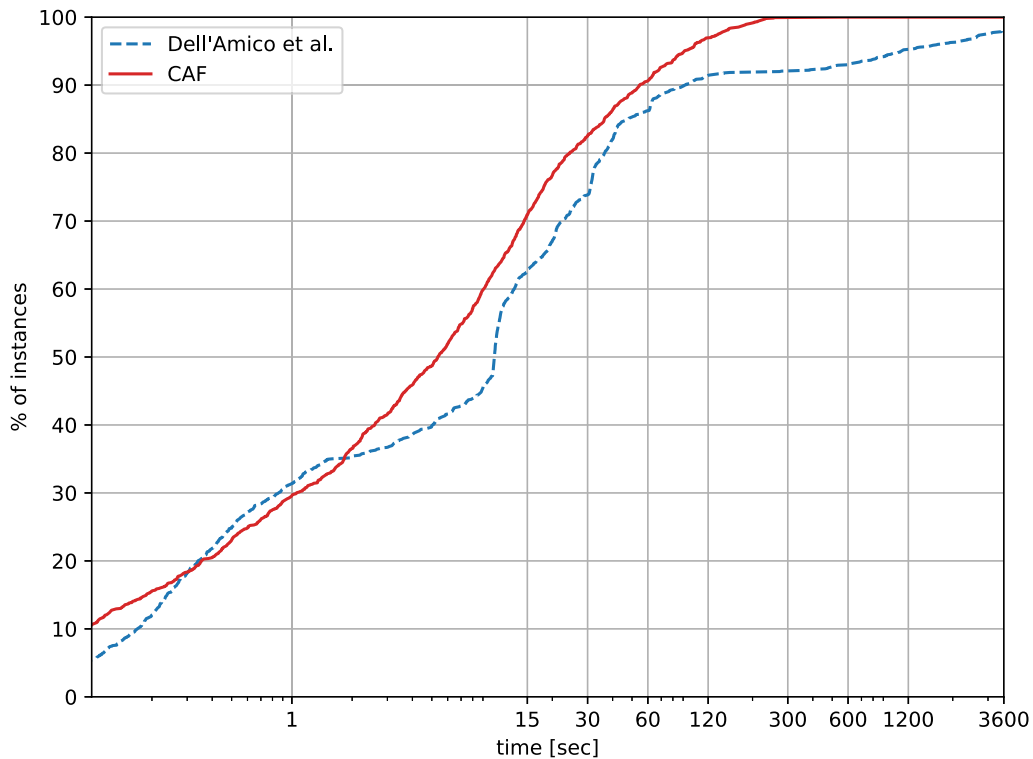


Fig. 8. Performance profile for the comparison of CAF and B&P+ from Dell'Amico et al. (2020) applied to Category (B).

$h := \lceil z_{LP}^{exp,*} \rceil$  for the number of servers in use to the model from Section 4 by requiring

$$\sum_{e \in \mathcal{E}^{out}(0, \emptyset)} \xi_e \geq h \quad \text{and} \quad \sum_{e \in \mathcal{E}^{FU}} \xi_e \geq h. \quad (22)$$

Observe that this bound can be obtained in relatively short time by solving the LP relaxation of the ordinary TBPP, and therefore it already appeared in any of the compact models on the TBPP-FU proposed in the literature<sup>2</sup>. Despite numerous improvements of these compact formulations, only about two thirds of the 160 classical test instances from Category (A1) could be solved to proven optimality so far. Currently, the most successful approach is a model of M1-type, see Section 4, whose final variant was recently presented in Martinovic & Strasdat (2022). The results listed in Table 5 contain the following main information:

- CAF can solve any single benchmark instance from Category (A1). Moreover, the solution times are (almost) always much better than reported in Martinovic & Strasdat (2022). Having a look at the overall average, we see that the solution time has reduced by more than 90% again.
- Since we have so far collected only exemplary model generation times for instances of Category (B), Table 5 also contains the respective average values  $t_{mod}$  for the instances considered here. It is again noticeable that these generation times are much smaller than the available time limit (and, thus, acceptable). Except for some very easy instance classes with a rather small number of items, where the time required for both formulations is in the very low seconds range, CAF wins the com-

<sup>2</sup> To be more precise, the lower bound on the number of servers (that is, the first inequality from (22)) directly appears in any compact formulation from the literature, whereas the lower bound on the number of fire-ups (that is, the second inequality from (22)) is implicitly imposed by means of valid inequalities, see Martinovic, Strasdat, & Selch (2021). The important overall message, then, is that CAF has no advantage based on modeling by adding the lower bounds presented in (22).

parison with the compact formulation (generally clearly) even when the modeling times are taken into account.

- For some of the  $(d_s, c_H)$  scenarios, we see that  $t_{mod} \approx t$  or even  $t_{mod} > t$  holds for CAF. However, as can be observed from the comparison with  $t_{mod}$  of other parameter configurations, this is not an indication of a disproportionately large modeling time. In fact, in these scenarios, the graph has many cliques of rather small cardinality (since the job durations are short), together with relatively few feasible patterns (since many items are incompatible) and, consequently, a somewhat small number of arcs. Thus, the optimization problems obtained for CAF generally have the fewest variables and can therefore be solved particularly fast. We will come back to this point with more details later in Table 8.
- In particular, those instances where there are many possible item combinations (i.e., the constellation  $(d_L, c_L)$ ) proved to be very difficult or even intractable for all the compact formulations, see Martinovic et al. (2022). In contrast, CAF can also solve the hardest subsets of these instances to proven optimality in less than ten minutes on average. However, we clearly see, especially when compared to the computation times of the ordinary TBPP reported in Table 3, that these are indeed already somewhat more challenging instances even for CAF. This is because, with such a large number of possible temporal interactions between the items, the maximal cliques generally consist of many different items and, thus, the number of nodes (as well as, implicitly, the number of arcs) increases considerably compared to other parameter constellations. Therefore, the modeling time is by far the largest for these instance groups, and the relatively large flow-based ILP to be solved then also requires some noticeable computation times.

Typically, another advantage of flow models over compact formulations is given by a substantially better LP bound, see de Lima et al. (2022a). Here, too, for Category (A1) a corresponding dominance relation can be manifested empirically in Table 6. However,

**Table 5**  
 Numerical comparison (for the TBPP-FU) between CAF and the best solution approach from the literature, that is, a compact M1-type model from Martinovic & Strasdat (2022), for Category (A1). In addition to the pure solution times, we also report about  $t_{mod}$ , the time to build the network and the corresponding ILP formulation. Moreover, the average clique size  $|C_i|_{avg}$  is given.

n	$\bar{s}$	$d_i$	$c_i$	$ C_i _{avg}$	CAF			Martinovic & Strasdat (2022)			
					$t_{mod}$	t	opt	t	opt		
50	50	$d_S$	$c_L$	21.1	1.0	1.6	(5)	4.0	(5)		
			$c_H$	18.7	0.2	0.1	(5)	1.3	(5)		
		$d_L$	$c_L$	33.6	1.7	0.8	(5)	360.8	(4)		
			$c_H$	34.9	0.5	0.1	(5)	0.5	(5)		
		60	$d_S$	$c_L$	16.8	0.6	0.9	(5)	1.2	(5)	
				$c_H$	17.4	0.3	0.2	(5)	1.7	(5)	
	$d_L$		$c_L$	33.2	2.0	1.2	(5)	11.2	(5)		
			$c_H$	30.9	0.5	0.1	(5)	0.6	(5)		
	Average (Sum)				25.8	0.9	<b>0.6</b>	<b>(40)</b>	47.7	(39)	
	100	100	$d_S$	$c_L$	20.0	2.1	4.2	(5)	3.6	(5)	
				$c_H$	20.1	0.7	0.3	(5)	78.6	(5)	
			$d_L$	$c_L$	38.6	10.9	78.2	(5)	1449.8	(1)	
$c_H$				39.1	3.0	1.7	(5)	1092.9	(2)		
120			$d_S$	$c_L$	16.9	1.3	3.1	(5)	85.3	(5)	
				$c_H$	16.6	0.5	0.3	(5)	83.6	(5)	
		$d_L$	$c_L$	31.1	5.6	13.7	(5)	685.9	(4)		
			$c_H$	32.7	1.7	0.7	(5)	546.9	(4)		
Average (Sum)				26.9	3.2	<b>12.8</b>	<b>(40)</b>	503.3	(31)		
150		150	$d_S$	$c_L$	19.5	3.2	13.2	(5)	85.7	(5)	
				$c_H$	20.1	1.0	0.6	(5)	1464.1	(1)	
			$d_L$	$c_L$	39.6	18.9	177.1	(5)	1462.2	(1)	
				$c_H$	38.5	4.2	4.2	(5)	1372.0	(2)	
			180	$d_S$	$c_L$	16.7	2.1	8.0	(5)	37.2	(5)
					$c_H$	17.0	0.9	0.4	(5)	853.9	(3)
		$d_L$		$c_L$	33.4	13.0	175.1	(5)	1198.4	(3)	
				$c_H$	32.9	3.0	2.2	(5)	1494.0	(1)	
		Average (Sum)				27.2	5.8	<b>47.6</b>	<b>(40)</b>	995.9	(21)
	200	200	$d_S$	$c_L$	20.4	5.1	47.1	(5)	99.2	(5)	
				$c_H$	19.8	1.5	1.6	(5)	1800.0	(0)	
			$d_L$	$c_L$	40.3	29.9	571.0	(5)	1800.0	(0)	
				$c_H$	39.6	7.3	9.1	(5)	1624.6	(1)	
			240	$d_S$	$c_L$	17.0	3.3	14.1	(5)	123.1	(5)
					$c_H$	16.9	1.2	0.7	(5)	1693.9	(2)
		$d_L$		$c_L$	33.2	18.6	283.2	(5)	1201.9	(2)	
				$c_H$	33.0	5.0	14.2	(5)	1800.0	(0)	
		Average (Sum)				27.5	9.0	<b>117.6</b>	<b>(40)</b>	1267.8	(15)
Total: Average (Sum)				26.9	4.7	<b>44.7</b>	<b>(160)</b>	703.7	(106)		

**Table 6**  
 Average rounded-up LP bound for instances of Category (A1), averaged over the input parameter ‘time horizon’. Hence, in this table, every number is the average of ten instances.

$d_i$	$c_i$	n = 50		n = 100		n = 150		n = 200	
		CAF	Martinovic & Strasdat (2022)	CAF	Martinovic & Strasdat (2022)	CAF	Martinovic & Strasdat (2022)	CAF	Martinovic & Strasdat (2022)
$d_S$	$c_L$	18.2	18.2	21.2	21.2	20.6	20.6	23.0	23.0
	$c_H$	24.9	24.6	32.7	31.0	36.7	34.6	39.0	35.1
$d_L$	$c_L$	29.8	29.8	32.6	32.6	35.2	35.2	35.2	35.2
	$c_H$	43.8	43.8	47.5	46.9	51.5	50.2	50.9	49.9
Average		<b>29.2</b>	29.1	<b>33.5</b>	32.9	<b>36.0</b>	35.2	<b>37.0</b>	35.8

except for the constellation ( $d_S, c_H$ ) being the most favorable setup in terms of a small network, the deviations are typically less than 3% and, therefore, not as considerable as one might have expected. This is due to the fact that the compact formulations have been significantly improved by numerous techniques like valid inequalities (Martinovic, Strasdat, & Selch, 2021; Martinovic et al., 2022). In particular, any ILP model from the literature already uses the lower bound of the exponential-size TBPP formulation, which is very powerful especially in those scenarios where an optimal solution contains a small number of servers and only a few additional fire-ups (like in the setting ( $d_S, c_L$ ), see also Table 5).

**Remark 10.** As a consequence of the very powerful model improvements (for the compact formulations from the literature) discussed above, we note that the LP bound of CAF does not dominate the bound of the compact formulation from Martinovic & Strasdat (2022) for any possible instance, so a general theoretical result cannot be established.

As a summary of the discussion of Category (A1) and in the light of the graphical illustration of computational results chosen in Martinovic, Strasdat, & Selch, 2021, Martinovic et al. (2022), we would also like to provide the following performance profile, in

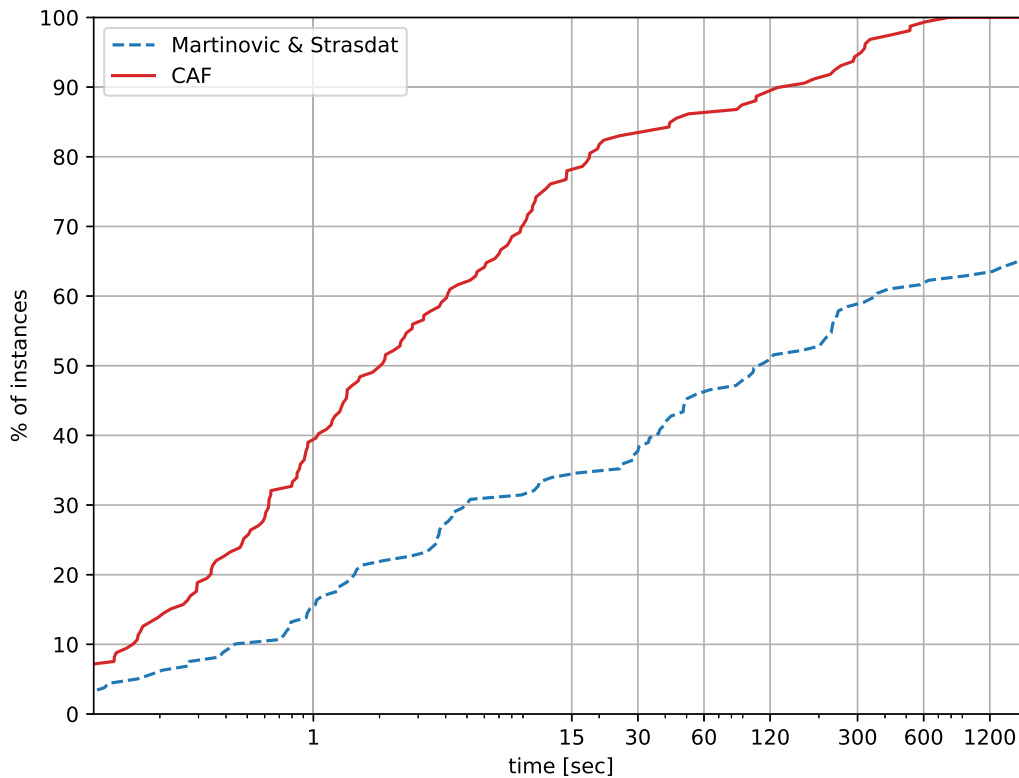


Fig. 9. Performance profile for the comparison of CAF and the compact model from Martinovic & Strasdat (2022) applied to Category (A1).

Table 7

Numerical results for instances of Category (B), always with  $t_{max} = 1800s$ . For the sake of an easier comparison, we repeat the solution times obtained for the classical TBPP from Table 3 in the first two columns of the main table. Moreover, we also display a selection of the results which would have been obtained with the approach from Martinovic & Strasdat (2022), but also point out that these results did not appear in the literature before. Due to this reason, we just conducted and included a reasonable subset of these additional calculations which is, however, sufficient to anticipate the general trends.

$ T_5^{nd} $	TBPP		TBPP-FU		Martinovic & Strasdat (2022)	
	t	opt	t	opt	t	opt
10	0.4	(100)	0.5	(100)	20.0	(99)
20	1.4	(100)	1.5	(100)	144.2	(94)
30	2.5	(100)	3.9	(100)	293.0	(90)
40	3.7	(100)	12.3	(100)	461.1	(82)
50	6.9	(100)	11.4	(100)	572.4	(75)
60	10.0	(100)	17.1	(100)	684.5	(67)
70	15.7	(100)	26.8	(100)	-	-
80	16.5	(100)	46.8	(100)	-	-
90	20.4	(100)	42.8	(100)	-	-
100	23.9	(100)	53.1	(100)	-	-
110	26.7	(100)	59.4	(100)	-	-
120	30.7	(100)	70.9	(100)	-	-
130	38.6	(100)	83.9	(100)	-	-
140	40.2	(100)	97.7	(100)	-	-
150	51.5	(100)	106.7	(100)	-	-
Average (Sum)	19.3	(1500)	42.3	(1500)	-	-

addition to the values appearing in Table 5. In Fig. 9, it can be seen that CAF is clearly ahead of the best compact formulation during the entire observation period. Remarkably, after only 15 seconds almost 80 % of all instances are already solved, also making CAF suitable for applications where decisions have to be made within rather short time.

In a last experiment, we interpret the instances of Category (B) as TBPP-FU instances, to enrich the variety of test sets for the latter problem, and collect the numerical results in Table 7. Although there are no corresponding calculations for the compact model from Martinovic & Strasdat (2022) reported in the

literature, we added some results for moderate instance sizes to enable a rough comparison. As can be seen in Table 7, neither of these instances is challenging for the TBPP-FU when using CAF, because (on average) they all can be solved in less than one minute, while even the harder subsets just require roughly the double amount of time. In particular, already for the 600 representative instances considered, CAF clearly outperforms the results which can be obtained for the compact model from Martinovic & Strasdat (2022), so that spending more computation time to fill in the associated columns completely is definitely not required.



**Table 8**

Average numbers of states and transitions (in units of  $10^3$ ) for the instances of Category (A). The averages are calculated based on ten instances each (since the criterion ‘time horizon’ is not specified here to keep the list short).

	$d_i$	$c_i$	$n = 50$	$n = 100$	$n = 150$	$n = 200$	$n = 500$	$n = 1000$
$ \mathcal{V} $	$d_S$	$c_L$	5.7	11.4	17.5	27.5	67.2	140.7
		$c_H$	1.6	4.0	6.8	9.7	25.5	50.5
	$d_L$	$c_L$	13.1	59.3	115.5	175.3	469.4	904.9
		$c_H$	3.4	16.9	25.9	44.1	122.6	240.1
$ \mathcal{E} $	$d_S$	$c_L$	7.7	15.8	24.2	37.8	91.4	191.5
		$c_H$	2.2	5.3	8.8	12.7	33.4	65.7
	$d_L$	$c_L$	17.6	71.4	137.3	206.5	553.1	1064.6
		$c_H$	4.4	20.1	30.2	51.4	141.5	276.4

An interesting side aspect of the results presented in Table 7 is that, for any  $|T_S^{nd}|$ , solving the TBPP-FU is more time-consuming than coping with the traditional TBPP. Since neither the associated networks nor the sets of variables and constraints differ very much (among the TBPP and the TBPP-FU), we partly attribute this to the observation that, as discussed earlier in Section 4, the TBPP-FU graph itself offers some symmetries to represent a feasible pattern. However, the more important reason is that any feasible schedule is now assessed by two terms (number of servers, number of fire-ups), typically entailing a wider variety of possible “numerical states” of the objective function, especially when there are much more fire-ups than active servers in an optimal solution.

6.3. Exploring the Limits: Numerical Results for Very Large Instances

In the previous subsection, we have demonstrated that CAF is able to solve all benchmark instances from the literature (that have been investigated so far in the context of exact approaches) in a short time, both for the TBPP and for the TBPP-FU. We therefore provide some more numerical tests, in particular to also show the limitations of our new approach:

- For Category (A), such instances have already been collected in the literature, see Aydin et al. (2020), but they have only been treated heuristically so far. To this end, these 80 instances already appeared in the presentation at the beginning of Section 4 as Category (A2).
- For Category (B), the relevant literature does not yet specify larger instances. Nevertheless, such instances can easily be obtained according to the known construction principles even for values  $|T_S^{nd}| > 150$ . To this end, applying the procedure described in Dell’Amico et al. (2020) to the raw data from Caprara et al. (2013), we obtain more instances along the lines of Category (B). To be more precise, we thus generated 100 instances each for any parameter value  $|T_S^{nd}| \in \{160, 170, \dots, 400\}$ . Although there is no set (B1), we will refer to these 2500 very large instances as Category (B2) to synchronize with the terminology of Category (A).

Let us proceed in alphabetical order. Before examining the actual numerical results of Category (A2), we would first like to present some data on the size of the respective CAF networks. For this purpose, we have summarized the number of nodes and arcs in Table 8, but we included all instances from Category (A) to provide a better overview of the overall evolution.

In particular, the following interesting insights should be noted:

- We see that the different parameter constellations lead to very heterogeneous graph sizes and thus cover a reasonable range of different benchmark scenarios. In particular, it becomes clear that the configuration  $(d_S, c_H)$  indeed leads to the smallest networks, while  $(d_L, c_L)$  typically allows for very many item interactions and thus requires many states and transitions. This

**Table 9**

Number of optimally solved instances and average computation times for Category (A2) for two different time limit settings: the classical  $t_{max} = 1800$  seconds vs. an open-end calculation with no time limit (indicated by  $t_{max} = \infty$ ).

$n$	$\bar{s}$	$d_i$	$c_i$	$t_{mod}$	$t_{max} = 1800s$		$t_{max} = \infty$	
					t	opt	t	opt
500	500	$d_S$	$c_L$	11.9	121.7	(5)	121.7	(5)
			$c_H$	4.6	4.7	(5)	4.7	(5)
		$d_L$	$c_L$	75.4	1696.3	(1)	5103.3	(5)
			$c_H$	19.6	144.2	(5)	144.2	(5)
	600	$d_S$	$c_L$	7.8	47.4	(5)	47.4	(5)
			$c_H$	2.6	2.3	(5)	2.3	(5)
		$d_L$	$c_L$	53.8	1467.3	(2)	2011.9	(5)
			$c_H$	13.8	47.0	(5)	47.0	(5)
Average (Sum)				23.7	441.4	(33)	935.3	(40)
1000	1000	$d_S$	$c_L$	24.5	887.5	(4)	933.4	(5)
			$c_H$	7.2	12.4	(5)	12.4	(5)
		$d_L$	$c_L$	150.0	1800.0	(0)	20483.1	(5)
			$c_H$	38.9	176.8	(5)	176.8	(5)
	1200	$d_S$	$c_L$	18.2	340.1	(5)	340.1	(5)
			$c_H$	5.5	8.3	(5)	8.3	(5)
		$d_L$	$c_L$	103.1	1800.0	(0)	5615.1	(5)
			$c_H$	28.1	126.6	(5)	126.6	(5)
Average (Sum)				46.9	644.0	(29)	3462.0	(40)
Total: Average (Sum)				35.3	542.7	(62)	2198.6	(80)

is consistent with the associated observations of modeling and solution times (see Table 5), which have already been partially addressed in the related discussion before.

- Compared to Category (B), we see that the instances from Category (A) can be judged as more challenging on average. To illustrate this more thoroughly, we consider the case of  $n = 500$  items as an example. Then the data from Table 8 prove that an associated CAF graph in the case of Category (A) has on average about 171 thousand nodes and 205 thousand arcs. If we now search Table 1 for instances from Category (B) with comparable or even larger item numbers (i.e., for example,  $|T_S^{nd}| = 100$  and Classes VI or X), we notice that the associated networks are generally much smaller. It is therefore to be expected that CAF will faster reach its limits in the case of Category (A), in particular since already for  $n = 1000$  items partly more than one million integer variables have to be dealt with according to Table 8.

In view of these remarks, a time limit of 30 minutes is no longer sufficient in some cases to cope with these very large instances. We have therefore also performed a calculation without any time limit in Table 9 to determine the so far unknown optimal value of these instances on the one hand, and to get a more precise impression of how long it actually takes to solve such challenging instances on the other hand. We highlight the following main observations:

**Table 10**  
Average size of the CAF network for some instances from Category (B2) depending on the class index.

$ T_S^{nd} $	number of nodes ( $ \mathcal{V} $ )				number of arcs ( $ \mathcal{E} $ )					
	200	250	300	350	400	200	250	300	350	400
I	6.5	8.4	10.0	11.8	13.6	8.0	10.3	12.2	14.5	16.7
II	19.6	24.1	27.6	31.7	35.8	26.5	32.4	37.2	42.7	48.2
III	34.3	44.6	55.6	62.8	71.7	44.8	58.7	73.0	82.5	94.4
IV	69.8	91.7	108.4	126.9	146.1	95.1	125.3	147.6	172.6	198.7
V	114.1	142.4	173.3	205.2	240.6	156.0	195.2	236.8	280.1	328.4
VI	83.8	107.1	129.6	151.3	171.2	179.2	229.7	277.5	324.9	367.4
VII	117.8	151.6	175.6	209.9	243.6	160.1	204.9	237.2	283.7	330.2
VIII	119.2	146.0	178.6	208.6	231.7	187.7	230.7	283.5	329.0	364.7
IX	82.4	104.4	123.6	145.3	165.2	184.0	234.0	275.5	322.5	366.8
X	194.4	248.3	299.2	343.1	386.6	301.6	386.1	467.3	535.6	602.1

- The modeling times of these huge instances are still perfectly fine, as even the most challenging subset requires only 2.5 minutes on average. This is a justifiable effort in view of the expected solution time and gives hope that, in the future, possibly also these instances can be coped with more efficiently as a consequence of the steady progress in terms of optimization software.
- Still, 62 of the 80 instances are solved optimally within a maximum of 30 minutes. The unsuccessfully attempted instances, with only one exception, all originate from the constellation  $(d_L, c_L)$ , which has already been identified as the greatest challenge before, with the help of Table 8. Comparing the two columns ( $t_{max} = 1800s$  vs.  $t_{max} = \infty$ ), it is noticeable that the only unsolved instance that did not come from that subset required a solution time that was roughly in the range of half an hour, so that one could also speak of a random effect here.
- Overall, it can be stated that on average all 80 instances are solved in less than one hour. The longest solution time we observed was slightly less than nine hours (for an instance with  $n = 1000$  items and the combination  $(d_L, c_L)$ ) – a time that is admittedly already relatively long, but nevertheless would not be achievable at all with the compact formulations from the literature.

Finally, we consider the instances from Category (B2) and first present an overview of the size of the corresponding CAF graphs in Table 10.

On the one hand, it is noticeable that these networks are on average much smaller than it was the case for the most difficult instances from Category (A2). By way of example, we mention that even for the most challenging combinations (e.g., Class X and  $|T_S^{nd}| = 400$ ), one detects a much smaller size on average than, say, for  $n = 1000$  items and the setting  $d_L$  in Category (A2). The instances from Category (B2) can thus still be classified as easier, despite a significant increase of the input parameter  $|T_S^{nd}|$ . It is also interesting to note that in contrast to Table 8, here a doubling of the input parameter  $|T_S^{nd}|$  also leads to a doubling of the number of nodes and arcs for basically any instance class. This is mainly due to the fact that the number of items per clique for each class is predefined by the  $a$ -parameter choice (see Dell’Amico et al., 2020 for the construction details) and thus, with increasing  $|T_S^{nd}|$ , only the number of cliques grows. Consequently, the graph size increases only proportionally to that value and not in a “combinatorial sense” (i.e., in the size of the cliques).

As a result of these observations, it seems reasonable to assume that even these enlarged instances will not be too challenging for CAF yet, and indeed this is also visible in the results depicted in Table 11. We highlight that in the case of TBPP, still every single instance can be solved optimally within a relatively short time (less than three minutes on average for any  $|T_S^{nd}|$ ). Conversely, within the time limit of 1800 seconds, a very few instances of the TBPP-

**Table 11**  
Numerical Results for instances of Category (B2) for both, the TBPP and the TBPP-FU (with  $t_{max} = 1800s$ ).

$ T_S^{nd} $	TBPP		TBPP-FU	
	t	opt	t	opt
160	48.9	(100)	123.1	(100)
170	55.7	(100)	133.7	(100)
180	54.2	(100)	139.6	(100)
190	62.8	(100)	158.6	(100)
200	71.3	(100)	201.2	(99)
210	80.9	(100)	183.0	(99)
220	73.6	(100)	205.0	(99)
230	85.6	(100)	206.4	(99)
240	82.4	(100)	226.2	(99)
250	100.6	(100)	241.8	(98)
260	97.9	(100)	267.5	(98)
270	99.6	(100)	266.6	(97)
280	106.7	(100)	314.2	(96)
290	109.6	(100)	287.3	(96)
300	115.2	(100)	310.1	(96)
310	120.0	(100)	327.3	(96)
320	119.3	(100)	363.3	(94)
330	131.2	(100)	329.6	(97)
340	136.8	(100)	338.7	(98)
350	135.3	(100)	376.0	(95)
360	138.9	(100)	363.5	(96)
370	140.3	(100)	343.9	(98)
380	161.5	(100)	425.4	(94)
390	166.8	(100)	415.8	(95)
400	164.4	(100)	404.4	(97)
Average (Sum)	106.4	(2500)	278.1	(2436)

FU can no longer be dealt with because, as described earlier, solving this problem is generally somewhat more costly. We would like to note that most of the unsolved instances are from Class X and therefore, as seen previously in Table 10, correspond to the (on average) largest ILP models. In addition, especially for larger parameter settings (that is, approximately,  $|T_S^{nd}| \geq 300$ ), sometimes one or two instances from Classes IV, V, or VII cannot be tackled successfully due to some random effects, so that there is no strict monotonicity in the number of instances solved to proven optimality for the TBPP-FU. However, and this is the difference to the instances from Category (A2), the maximum computation time in our case was only about 142 minutes (for one instance from Class X), so that we can assume that CAF will also solve the vast majority of even larger benchmark instances constructed according to the same principles.

**Remark 11.** To get a somewhat more accurate idea of the actual boundaries of applying CAF to Category (B2), we conducted some further (less systematic) internal tests. In these calculations, we observed some first memory issues (in terms of storing the resulting branch-and-bound trees) when dealing with TBPP-FU instances

having  $|T_S^{nd}| = 500$  non-dominated starting times. However, this only happened in a very few exceptional cases, so that even here the size of the graph is not problematic, in general, and almost all instances can still be tackled properly.

### 7. Conclusions

In this article, we addressed the exact solution of two types of temporal bin packing problems, the TBPP and the TBPP-FU, by developing a new graph-theoretic approach (called CAF). Such an idea had previously been identified in the literature as an inefficient solution method given the generally exponential size of the resulting networks. By cleverly grouping equivalent states together in the construction of the graph, we managed to significantly reduce the number of nodes and arcs compared to previous concepts from Caprara et al. (2013) and Clautiaux et al. (2021). Remarkably, the associated ILP formulations can now be generated in a relatively short time even for very large instances, and thus they can easily be passed to a commercial ILP solver. Based on extensive test calculations, it turns out that for the first time ever all benchmark instances of the TBPP and the TBPP-FU, previously used in the context of exact approaches, can be solved to proven optimality in reasonably short time. Moreover, our new formulation not only outperforms the previous state of the art in terms of solution times, but also succeeds in handling much larger new benchmark instances based on the classical test scenarios mentioned before. All in all, we have thus presented a powerful unified approach for solving temporal bin packing problems, the basic concepts of which can be prospectively applied (with minor modifications, if necessary) to other classes of optimization problems in the field of interval scheduling. In future research, we will try to further improve this very promising concept, for example by incorporating reduced cost variable fixing or by investigating whether the now known optimal solutions can also be obtained using thinned-out graphs, like for example illustrated for the CSP in de Lima et al. (2022b); Delorme & Iori (2020).

### Declarations

**Funding:** This work has been supported by FCT – Fundação para a Ciência e Tecnologia within the R&D Units Project Scope UIDB/00319/2020.

**Conflicts of interest:** The authors declare that they do not have any conflicts of interest.

**Availability of data and material:** For clarity, the instances used in this paper have been gathered together at <https://github.com/wotzlaff/tbpp-instances>. However, we note that most of these instances were originally designed in Aydin et al. (2020) and Dell’Amico et al. (2020), and some of them were already available online, see <https://github.com/sibirbil/TemporalBinPacking>.

**Code availability:** The instances were solved by the commercial software Gurobi. The underlying implementation of the models in Python can be found at <https://github.com/wotzlaff/tbpp-caf>.

### Appendix A. Further Illustrations

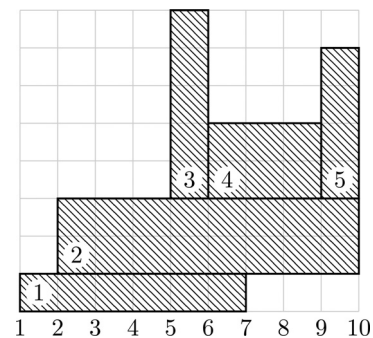


Fig. A.1. Visualization of the instance  $E_2$  from Example 5.

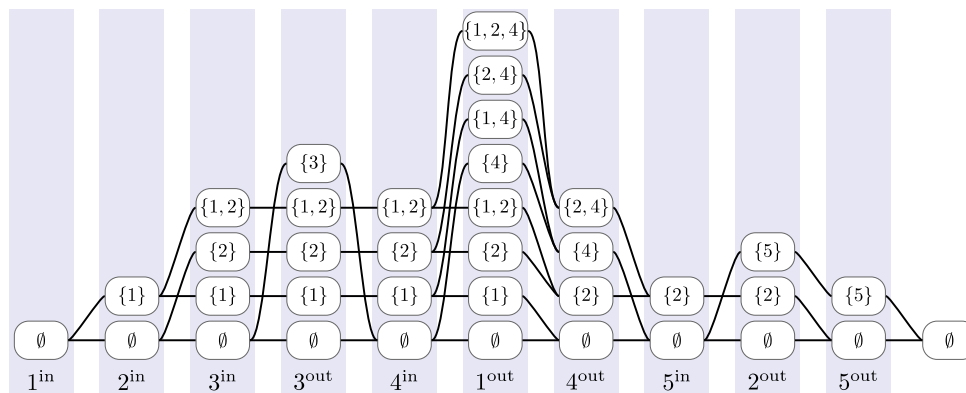


Fig. A.2. The network from Clautiaux et al. (2021) when applied to instance  $E_2$  from Example 5.

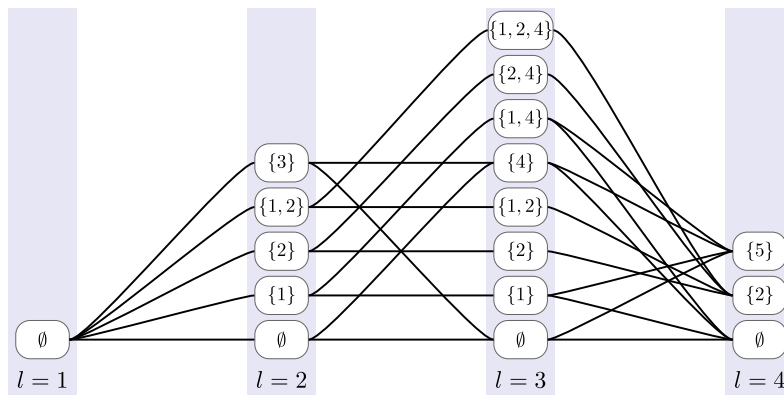


Fig. A.3. The network from Caprara et al. (2013) when applied to the instance  $E_2$  from Example 5.

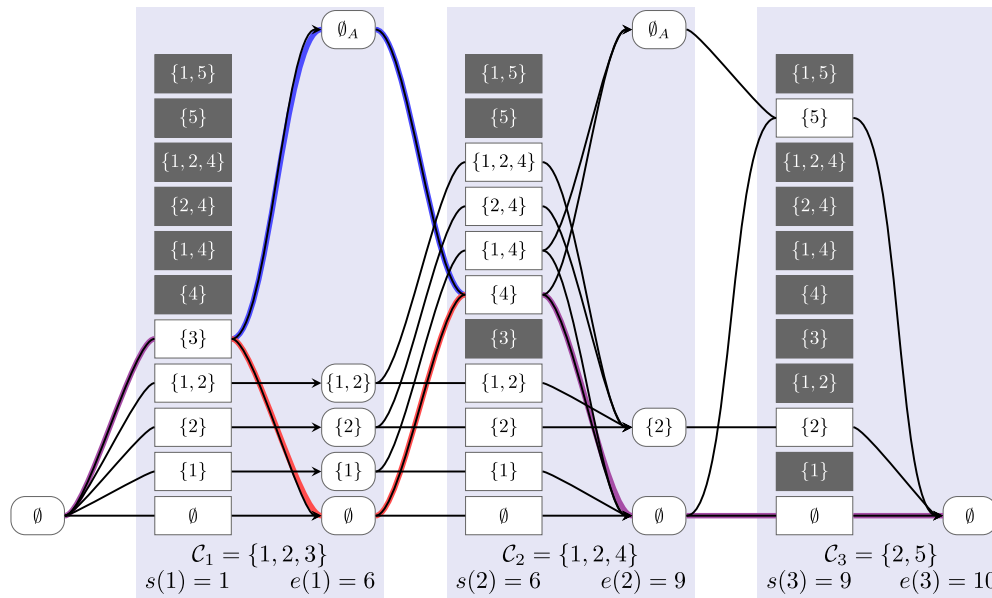


Fig. A.4. The combinatorial arcflow graph for the TBPP-FU and the instance  $E_2$  with two different paths leading to the same pattern  $J = \{3, 4\}$ . The joint (parts of the) arcs are colored purple, while the two alternatives to pass layer  $l = 1$  are painted red and blue, respectively. In terms of optimization, the blue path will be preferred, because it does not imply additional fire-up costs (which is consistent with the pattern  $J$  described before). (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

References

Angelelli, E., Bianchessi, N., & Filippi, C. (2014). Optimal interval scheduling with a resource constraint. *Computers & Operations Research*, 51, 268–281.

Arkin, E. M., & Silverberg, E. B. (1987). Scheduling jobs with fixed start and end times. *Discrete Applied Mathematics*, 18(1), 1–8.

Aydin, N., Muter, I., & Ilker Birbil, S. (2020). Multi-objective temporal bin packing problem: An application in cloud computing. *Computers & Operations Research*, 121.

Bar-Noy, A., Bar-Yehuda, R., Freund, A., & Naor, J. (2001). A unified approach to approximating resource allocation and scheduling. *Journal of the ACM*, 48(5), 1069–1090.

Bar-Noy, A., Canetti, R., Kutten, S., Mansour, Y., & Schieber, B. (1999). Bandwidth allocation with preemption. *SIAM Journal on Computing*, 28(5), 1806–1828.

Barnett, T., Jr., Jain, S., Andra, U., & Khurana, T. (2018). Cisco visual networking index (VNI) complete forecast update, 2017–2022. *APJC Cisco Knowledge Network (CKN) Presentation*. ([https://www.cisco.com/c/dam/m/en\\_us/network-intelligence/service-provider/digital-transformation/knowledge-network-webinars/pdfs/1213-business-services-ckn.pdf](https://www.cisco.com/c/dam/m/en_us/network-intelligence/service-provider/digital-transformation/knowledge-network-webinars/pdfs/1213-business-services-ckn.pdf))

Barnhart, C., Johnson, E. L., Nemhauser, G. L., Savelsbergh, M. W. P., & Vance, P. H. (1998). Branch-and-price: Column generation for solving huge integer programs. *Operations Research*, 46(3), 316–329.

Bartlett, M., Frisch, A. M., Hamadi, Y., Miguel, I., Tarim, S., & Unsworth, C. (2005). The temporal knapsack problem and its solution. *Lecture Notes in Computer Science*, 3524, 34–48.

Biedl, T. (2005). Graph-theoretic algorithms. In *Lecture notes*. University of Waterloo.

Boschetti, M. A., Hadjiconstantinou, E., & Mingozzi, A. (2002). New upper bounds

for the two-dimensional orthogonal non guillotine cutting stock problem. *IMA Journal of Management Mathematics*, 13(2), 95–119.

Buyya, R., Yeo, C. S., Venugopal, S., Broberg, J., & Brandic, I. (2009). Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation Computer Systems*, 25(6), 599–616.

Caprara, A., Furini, F., & Malaguti, E. (2013). Uncommon Dantzig-Wolfe reformulation for the temporal knapsack problem. *INFORMS Journal on Computing*, 25(3), 560–571.

Caprara, A., Furini, F., Malaguti, E., & Traversi, E. (2016). Solving the temporal knapsack problem via recursive Dantzig-Wolfe reformulation. *Information Processing Letters*, 116(5), 379–386.

Chen, B., Hassin, R., & Tzur, M. (2002). Allocation of bandwidth and storage. *IEE Transactions*, 24, 501–507.

Clausen, J. V., Lusby, R., & Ropke, S. (2022). Consistency cuts for Dantzig-Wolfe reformulations. To appear in: *Operations Research* (<https://doi.org/10.1287/opre.2021.2160>).

Clautiaux, F., Carlier, J., & Moukrim, A. (2007). New reduction procedures and lower bounds for the two-dimensional bin packing problem with fixed orientation. *Computers & Operations Research*, 34(8), 2223–2250.

Clautiaux, F., Detienne, B., & Guillot, G. (2021). An iterative dynamic programming approach for the temporal knapsack problem. *European Journal of Operational Research*, 293(2), 442–456.

Darmann, A., Pferschy, U., & Schauer, J. (2010). Resource allocation with time intervals. *Theoretical Computer Science*, 411, 4217–4234.

de Cauwer, M., Mehta, D., & O’Sullivan, B. (2016). The temporal bin packing problem: An application to workload management in data centres. In *Proceedings of the 28th IEEE international conference on tools with artificial intelligence* (pp. 157–164).

- de Lima, V. L., Alves, C., Clautiaux, F., Iori, M., & Valério de Carvalho, J. M. (2022a). Arc flow formulations based on dynamic programming: Theoretical foundations and applications. *European Journal of Operational Research*, 296(1), 3–21.
- de Lima, V. L., Iori, M., & Miyazawa, F. K. (2022b). Exact solution of network flow models with strong relaxations. To appear in: *Mathematical Programming* (<https://doi.org/10.1007/s10107-022-01785-9>).
- Dell'Amico, M., Delorme, M., Iori, M., & Martello, S. (2019). Mathematical models and decomposition methods for the multiple knapsack problem. *European Journal of Operational Research*, 274(3), 886–899.
- Dell'Amico, M., Furini, F., & Iori, M. (2020). A branch-and-price algorithm for the temporal bin packing problem. *Computers & Operations Research*, 114.
- Delorme, M., & Iori, M. (2020). Enhanced pseudo-polynomial formulations for bin packing and cutting stock problems. *INFORMS Journal on Computing*, 32(1), 101–119.
- Delorme, M., Iori, M., & Martello, S. (2016). Bin packing and cutting stock problems: Mathematical models and exact algorithms. *European Journal of Operational Research*, 255, 1–20.
- Fettweis, G., Dörpinghaus, M., Castrillon, J., Kumar, A., Baier, C., Bock, K., ... Strufe, T. (2019). Architecture and advanced electronics pathways towards highly adaptive energy-efficient computing. *Proceedings of the IEEE*, 107(1), 204–231.
- Furini, F. (2011). Decomposition and reformulation of integer linear programming problems. Phd thesis, Università di Bologna.
- Gilmore, P. C., & Gomory, R. E. (1961). A linear programming approach to the cutting-stock problem (Part 1). *Operations Research*, 9, 849–859.
- Gschwind, T., & Irnich, S. (2017). Stabilized column generation for the temporal knapsack problem using dual-optimal inequalities. *OR Spectrum*, 39, 541–556.
- Hall, N. G., & Magazine, M. J. (1994). Maximizing the value of a space mission. *European Journal of Operational Research*, 78, 224–241.
- Kantorovich, L. V. (1939). Mathematical methods of organising and planning production. *Management Science*, 6, 366–422. (Russian, 1960 English)
- Kaplan, J. M., Forrest, W., & Kindler, N. (2008). Revolutionizing data center energy efficiency. Technical report, McKinsey & Company.
- Kellerer, H., Pferschy, U., & Pisinger, D. (2004). *Knapsack problems*. Springer.
- Kolen, A. W. J., Lenstra, J. K., Papadimitriou, C. H., & Spiessma, F. C. R. (2007). Interval scheduling: A survey. *Naval Research Logistics*, 54(5), 530–543.
- Martinovic, J., Delorme, M., Iori, M., Scheithauer, G., & Strasdat, N. (2020). Improved flow-based formulations for the skiving stock problem. *Computers & Operations Research*, 113.
- Martinovic, J., & Strasdat, N. (2022). Theoretical insights and a new class of valid inequalities for the temporal bin packing problem with fire-ups. *Preprint MATH-NM-01-2022, Technische Universität Dresden*. ([http://www.optimization-online.org/DB\\_HTML/2022/02/8791.html](http://www.optimization-online.org/DB_HTML/2022/02/8791.html)).
- Martinovic, J., Strasdat, N., & Selch, M. (2021). Compact integer linear programming formulations for the temporal bin packing problem with fire-ups. *Computers & Operations Research*, 132.
- Martinovic, J., Strasdat, N., Valério de Carvalho, J. M., & Furini, F. (2022). Variable and constraint reduction techniques for the temporal bin packing problem with fire-ups. *Optimization Letters*, 16, 2333–2358.
- Phillips, C. A., Uma, R. N., & Wein, J. (2000). Offline admission control for general scheduling problems. *Journal of Scheduling*, 3, 365–381.
- Ryan, D. M., & Foster, B. A. (1981). An integer programming approach to scheduling. In *Computer scheduling of public transportation urban passenger vehicle and crew scheduling* (pp. 269–280).
- Sadykov, R., Vanderbeck, F., Pessoa, A., Tahiri, I., & Uchoa, E. (2019). Primal heuristics for branch and price: The assets of diving methods. *INFORMS Journal on Computing*, 31(2), 251–267.
- Scheithauer, G. (2018). Introduction to cutting and packing optimization – problems, modeling approaches, solution methods. In *International series in operations research & management science* 263 (1st ed.), Springer.
- Valério de Carvalho, J. M. (2002). LP models for bin packing and cutting stock problems. *European Journal of Operations Research*, 141(2), 253–273.