



Optimization framework for DFG-based automated process discovery approaches

Adriano Augusto^{1,2} · Marlon Dumas² · Marcello La Rosa¹ · Sander J. J. Leemans³ ·
Seppe K. L. M. vanden Broucke^{4,5}

Received: 24 February 2020 / Revised: 27 September 2020 / Accepted: 18 November 2020 / Published online: 27 February 2021
© The Author(s) 2021

Abstract

The problem of automatically discovering business process models from event logs has been intensely investigated in the past two decades, leading to a wide range of approaches that strike various trade-offs between accuracy, model complexity, and execution time. A few studies have suggested that the accuracy of automated process discovery approaches can be enhanced by means of metaheuristic optimization techniques. However, these studies have remained at the level of proposals without validation on real-life datasets or they have only considered one metaheuristic in isolation. This article presents a metaheuristic optimization framework for automated process discovery. The key idea of the framework is to construct a directly-follows graph (DFG) from the event log, to perturb this DFG so as to generate new candidate solutions, and to apply a DFG-based automated process discovery approach in order to derive a process model from each DFG. The framework can be instantiated by linking it to an automated process discovery approach, an optimization metaheuristic, and the quality measure to be optimized (e.g., fitness, precision, F-score). The article considers several instantiations of the framework corresponding to four optimization metaheuristics, three automated process discovery approaches (Inductive Miner—directly-follows, Fodina, and Split Miner), and one accuracy measure (Markovian F-score). These framework instances are compared using a set of 20 real-life event logs. The evaluation shows that metaheuristic optimization consistently yields visible improvements in F-score for all the three automated process discovery approaches, at the cost of execution times in the order of minutes, versus seconds for the baseline approaches.

Keywords Automated process discovery · Metaheuristic optimization · Process mining

1 Introduction

Modern information systems such as customer relationship management (CRM) and enterprise resource planning (ERP) systems record transactions corresponding to activities executed within the business processes that these systems support. For example, a CRM system typically records transactions corresponding to the creation of a customer lead, a

request for quote, and various other activities related to customer leads, quotes, and purchase orders. These transactional records can be extracted via SQL queries or via dedicated application programming interfaces (APIs) and used to analyze the execution of the business processes supported by the CRM system, such as the lead-to-quote or the quote-to-order process.

Process mining [33] is a family of techniques to analyze transactional records associated to a given business process, also known as an event log, in order to extract insights about the performance of the process. Among other things, process mining techniques allow us to discover a process model from an event log, an operation known as *automated process discovery*. Automatically discovered process models allow analysts to understand how the process is executed in reality and to uncover unexpected behavior. When enhanced with performance information (e.g., average activity durations or

Communicated by Ulrich Frank.

✉ Adriano Augusto
a.augusto@unimelb.edu.au

¹ University of Melbourne, Melbourne, Australia

² University of Tartu, Tartu, Estonia

³ Queensland University of Technology, Brisbane, Australia

⁴ UGent, Ghent, Belgium

⁵ KU Leuven, Leuven, Belgium

waiting times), such models are also used for performance analysis, e.g., to detect bottlenecks.

The problem of automated process discovery has been intensely studied in the past two decades [8]. Research in this field has led to a wide range of automated process discovery approaches (APDAs) that strike various trade-offs between accuracy,¹ model complexity, and execution time. Existing approaches in this field rely on parameters (with certain default values) to strike this tradeoff. Analysts need to fine-tune these parameters to find a model with the best possible trade-off between different model quality metrics. This article addresses the question of how to automate the fine-tuning of automated process discovery techniques.

A few studies have suggested that the accuracy of APDAs can be enhanced by applying optimization metaheuristics. Early studies in this direction considered population-based metaheuristics (P-metaheuristics), chiefly genetic algorithms [13,16]. However, these heuristics are computationally heavy, requiring execution times in the order of hours to converge when applied to real-life logs [8]. Such high execution times make these techniques inapplicable in the context of exploratory and interactive process discovery, where an analyst may need to discover process models corresponding to several variants of a process (e.g., one process model per type of product, per type of customer, or per region or country) in order to compare the behavior of the process in different settings. Accordingly, other studies have considered the use of single-solution-based metaheuristics (S-metaheuristics) such as simulated annealing [18,29], which are less computationally demanding. However, these latter studies remain at the level of proposals without validation on real-life logs and comparison of trade-offs between alternative metaheuristics.

In this setting, this article studies the following question: *to what extent can the accuracy of APDAs be improved by applying single-solution-based metaheuristics?* To address this question, the article outlines a framework to enhance APDAs by applying optimization metaheuristics. The core idea is to perturb the intermediate representation of event logs used by several of the available APDAs, namely the directly-follows graph (DFG). The paper specifically considers perturbations that add or remove edges with the aim of improving fitness or precision, and in a way that allows the underlying APDA to discover a process model from the perturbed DFG.

The proposed framework can be instantiated by linking it to three components: (i) an automated process discovery approach; (ii) an optimization metaheuristic; and (iii) the quality measure to be optimized, such as fitness, pre-

cision, or F-score. The article considers instantiations of the framework corresponding to three APDAs (Inductive Miner [24],² Fodina [34], and Split Miner [9]), four optimization metaheuristics (iterative local search, repeated local search, tabu search, simulated annealing), and one accuracy measure (Markovian F-score).

Using a benchmark of 20 real-life logs, the article compares the accuracy gains yielded by the above optimization metaheuristics relative to each other, and relative to the baseline (unoptimized) APDAs upon which they rely. The experimental evaluation also considers the impact of metaheuristic optimization on model complexity measures as well as on execution times.

This article is an extended and revised version of a conference paper [10]. In the conference paper, we presented an approach to optimize the accuracy of one automated process discovery approach, namely Split Miner, by applying S-metaheuristics, and we reported a comparison between the benefits of applying single-solution-based metaheuristics against the benefits of applying P-metaheuristics (using Evolutionary Tree Miner [13] as representative APDA of this category). Our former comparison [10] showed that S-metaheuristics outperform P-metaheuristics not only in terms of execution time efficiency, but also in terms of accuracy of the discovered process models; such a result also supported the findings of the latest literature review of automated process discovery approaches [8]. This article extends our previous approach [10] into a modular framework that can be used to optimize other APDAs, specifically those APDAs that construct a DFG from the event log and use it as an intermediate artifact to discover a process model. This article also extends the conference paper by considering not only Split Miner, but also two other APDAs, namely Fodina and Inductive Miner. Finally, the article reports an empirical evaluation covering all three approaches (Split Miner, Fodina, and Inductive Miner). The evaluation not only proves the applicability and relevance of S-metaheuristics to the problem of automated process discovery, but it also highlights the benefits yielded by S-metaheuristics.

The rest of the paper is structured as follows. The next section gives an overview of APDAs and optimization metaheuristics, where we discuss the background and the related work. Section 3 presents the proposed metaheuristic optimization framework and its instantiations. Section 4 reports on the empirical evaluation, and Sect. 5 draws conclusions and future work directions.

¹ Here, the term accuracy is used in its informal sense to refer to how well a given process model reflects the event log from which it was discovered. Later in the article, we introduce specific measures of accuracy such as fitness and precision.

² We consider a specific version of Inductive Miner known as “Inductive Miner—directly-follows”.

2 Background and related work

In this section, we give an overview of existing approaches for automated process discovery, followed by an introduction to optimization metaheuristics in general, and their application to automated process discovery in particular.

2.1 Automated process discovery

The execution of business processes is often recorded in the form of *event logs*. An event log is a collection of event records produced by individual instances (i.e., cases) of the process. The goal of automated process discovery is to generate a process model that captures the behavior observed in or implied by an *event log*. To assess the goodness of a discovered process model, four quality dimensions are used [33]: fitness, precision, generalization, and complexity. *Fitness* (a.k.a. recall) measures the amount of behavior observed in the log that is captured by the model. A perfectly fitting process model is one that recognizes every trace in the log. *Precision* measures the amount of behavior captured in the process model that is observed in the log. A perfectly precise model is one that recognizes only traces that are observed in the log. *Generalization* measures to what extent the process model captures behavior that, despite not being observed in the log, is implied by it. Finally, *complexity* measures the understandability of a process model, and it is typically measured via size and structural measures. In this paper, we focus on fitness, precision, and F-score (the harmonic mean of fitness and precision).

A recent comparison of state-of-the-art APDAs [8] showed that an approach capable of consistently discovering models with the best fitness-precision trade-off is currently missing. The same study showed, however, that we can obtain consistently good trade-offs by hyperparameter-optimizing some of the existing APDAs based on DFGs—Inductive Miner [24], Structured Heuristics Miner [7], Fodina [34], and Split Miner [9]. These algorithms have a hyperparameter to tune the amount of filtering applied when constructing the DFG. Optimizing this and other hyperparameters via greedy search [8], local search strategies [14], or sensitivity analysis techniques [27], can greatly improve the accuracy of the discovered process models. Accordingly, in the evaluation reported later we use a hyperparameter-optimized version of Split Miner as one of the baselines.

The problem of discovering accurate process models from event logs is inevitably related to that of ensuring event log quality. There is a rich collection of methods for detecting and handling data quality issues in event logs [?]. However, this latter body of work is largely orthogonal to the contribution of this article, as this article focuses on discovering process models assuming that data quality issues have been addressed. This having been said, the methods

presented in this paper partially address one type of data quality issue, namely the presence of noise (infrequent behavior) in an event log [?]. To mitigate the impact of noise on the discovered process model, automated process discovery approaches, including those extended in this paper, apply a dependency filtering step. The optimization techniques proposed in this article fine-tune the level of filtering in order to maximize the accuracy of the discovered process model.

2.2 Optimization metaheuristics

The term *optimization metaheuristics* refers to a parameterized algorithm, which can be instantiated to address a wide range of optimization problems. Metaheuristics are usually classified into two broad categories [12]: (i) *single-solution-based metaheuristics*, or S-metaheuristics, which explore the solution space one solution at a time starting from a single initial solution of the problem; and (ii) *population-based metaheuristics*, or P-metaheuristics, which explore a population of solutions generated by mutating, combining, and/or improving previously identified solutions. S-metaheuristics tend to converge faster toward an optimal solution (either local or global) than P-metaheuristics, since the latter by dealing with a set of solutions require more time to assess and improve the quality of each single solution. P-metaheuristics are computationally heavier than S-metaheuristics, but they are more likely to escape local optima. Providing an exhaustive discussion on all the available metaheuristics is beyond the scope of this paper, in the following, we focus on the four S-metaheuristics that we integrated in our optimization framework and on the P-metaheuristics that have been previously adapted to address the problem of automated process discovery.

Iterated Local Search [30] starts from a (random) solution and explores the neighboring solutions (i.e., solutions obtained by applying a change to the current solution) in search of a better one. When a better solution cannot be found, it perturbs the current solution and starts again. The perturbation is meant to avoid local optimal solutions. The exploration of the solution-space ends when a given termination criterion is met (e.g., maximum iterations, timeout).

Tabu Search [19] is a memory-driven local search. Its initialization includes a (random) solution and three memories: short, intermediate, and long term. The short-term memory keeps track of recent solutions and prohibits to revisit them. The intermediate-term memory contains criteria driving the search toward the best solutions. The long-term memory contains characteristics that have often been found in many visited solutions, to avoid revisiting similar solutions. Using these memories, the neighborhood of the initial solution is explored and a new solution is selected accordingly. The solution-space exploration is repeated until a termination criterion is met.

Simulated Annealing [22] is based on the concepts of temperature (T , a parameter chosen arbitrarily) and energy (E , the objective function to minimize). At each iteration the algorithm explores (some of) the neighboring solutions and compares their energies with the one of the current solution. This latter is updated if the energy of a neighbor is lower, or with a probability that is function of T and the energies of the current and candidate solutions, usually $e^{-\frac{|E_1-E_2|}{T}}$. The temperature drops over time, thus reducing the chance of updating the current solution with a higher-energy one. The algorithm ends when a termination criterion is met, which often relates to the energy or the temperature (e.g., energy below a threshold or $T = 0$).

Evolutionary (Genetic) Algorithms [11,20] are inspired by Darwin's theory of evolution. Starting from a set of (random) solutions, a new solution is generated by mixing characteristics of two parents selected from the set of the current solutions, such an operation is known as *crossover*. Subsequently, *mutations* are applied to the new solutions to introduce randomness and avoid local optimal solution. Finally, the solutions obtained are assessed and a subset is retained for the next iteration. The algorithm continues until a stop criterion is met.

Swarm Particle Optimization [21] starts from a set of (random) solutions, referred to as particles. Each particle is identified using the concepts of *position* and *velocity*. The position is a proxy for the particle qualities and it embeds the characteristics of the solution, while the velocity is used to alter the position of the particles at each iteration. Furthermore, each particle has memory of its best position encountered during the roaming of the search space, as well as the best position encountered by any other particle. At each iteration, the algorithm updates the particles positions according to their velocities and updates the best positions found. When a termination condition is met, the algorithm returns the particle having the absolute best position among the whole swarm.

Imperialist Competitive Algorithm [4] is inspired by the historical colonial period. It starts from a (random) set of solutions, called countries. Each country is assessed via an objective function, and a subset is selected as imperialistic countries (the selection is based on their objective function scores). All the countries left (i.e., those having low objective function scores) are considered colonies of the closest (by characteristics) imperialistic country. Then, each colony is altered to resemble its imperialistic country, the objective function scores are re-computed, and the colonies that became better than their imperialistic country are promoted to imperialistic countries and vice-versa. When a termination condition is met, the country with the highest objective function score is selected as the best solution.

2.3 Optimization metaheuristics in automated process discovery

Optimization metaheuristics have been considered in a few previous studies on automated process discovery. An early attempt to apply P-metaheuristics to automated process discovery was the Genetic Miner proposed by de Medeiros [16], subsequently overtaken by the Evolutionary Tree Miner [13]. Other applications of P-metaheuristics include the contribution of Alizadeh and Norani [3] who showed to improve fitness and precision of the discovered process models by using the imperialist competitive algorithm, outperforming some state-of-the-art APDAs (including $\alpha++$ [36], Flexible Heuristics Miner [35], and Fodina [34]); however, the implementation of the method designed by Alizadeh et al. is not publicly available, and the benchmark they used differ from the one suggested in the latest literature review [8]. Some research studies adapted the particle swarm optimization metaheuristic to solve the problem of automated process discovery from event logs [15,17], but these studies are seminal and they lack of a solid evaluation on real-life logs. One of the most recent studies tried to combine evolutionary computation with particle swarm optimization [25] by extending the work of Buijs et al. [13], but also in this case the authors did not provide a working implementation of their method, and they did not assess it on public datasets, so that it is difficult to estimate the real benefits of their proposed improvements. In our context, the main limitation of P-metaheuristics is that they are computationally heavy due to the cost of constructing a solution (i.e., a process model) and evaluating its accuracy. This leads to execution times in the order of hours, to converge to a solution that in the end is comparable to those obtained by state-of-the-art APDAs that do not rely on optimization metaheuristics [8].

Finally, a handful of studies have considered the use of S-metaheuristics to automatically discover optimal process models, specifically simulated annealing [18,29], but these proposals are preliminary and have not been compared against state-of-the-art approaches on real-life logs.

3 Metaheuristic optimization framework

This section outlines our framework for optimizing APDAs by means of S-metaheuristics (cf. Sect. 2). First, we give an overview of the framework and its core components. Next, we discuss the adaptation of the S-metaheuristics to the problem of process discovery. Finally, we describe the instantiations of our framework for Split Miner, Fodina, and Inductive Miner.

3.1 Preliminaries

In order to discover a process model, an APDA takes as input an event log and transforms it into an intermediate representation from which a process model is derived. Below, we define one of the most popular intermediate representations, that is the directly-follows graph (DFG). Although other intermediate representations are available in the literature (e.g., behavioral profiles [28]), our framework focuses only on DFGs for two main reasons: first, because they are adopted by many state-of-the-art automated process discovery approaches [7,9,24,34,35]; second, because they allow us to leverage the Markovian accuracy [5] to facilitate the application of metaheuristics and the navigation of the solution space as we show later in this section.

Definition 1 (Event Log) Given a set of activities \mathcal{A} , an event log \mathcal{L} is a multiset of traces where a trace $t \in \mathcal{L}$ is a sequence of activities $t = \langle a_1, a_2, \dots, a_n \rangle$, with $a_i \in \mathcal{A}$, $1 \leq i \leq n$.

Definition 2 [Directly-follows graph (DFG)] Given an event log \mathcal{L} , its *directly-follows graph (DFG)* is a directed graph $\mathcal{G} = (N, E)$, where: N is the set of nodes, $N = \{a \in \mathcal{A} \mid \exists t \in \mathcal{L} \wedge a \in t\}$; and E is the set of edges $E = \{(x, y) \in N \times N \mid \exists t = \langle a_1, a_2, \dots, a_n \rangle, t \in \mathcal{L} \wedge a_i = x \wedge a_{i+1} = y [1 \leq i \leq n - 1]\}$.

By definition, each node of the DFG represents an activity recorded in at least one trace of the event log, while each edge of a DFG represents a directly-follows relation between two activities (represented by the node source and the node target of the edge). An APDA is said to be *DFG-based* if it first generates the DFG of the event log, then applies an algorithm to manipulate the DFG (e.g., removing edges), and finally converts the processed DFG into a process model. Such a processed DFG will not adhere any more to Definition 2; therefore, we redefine it as *Refined DFG*.

Definition 3 (Refined DFG) Given an event log \mathcal{L} and its DFG $\mathcal{G}_{\mathcal{L}} = (N, E)$, a *Refined DFG* is a directed graph $\mathcal{G} = (N', E')$, where: $N' \subseteq N$ and $E' \subseteq E$. If $N' = N$ and $E' = E$, the refined DFG is equivalent to the event log DFG.

Examples of DFG-based APDAs are Inductive Miner [24], Heuristics Miner [7,35], Fodina [34], and Split Miner [9]. Different DFG-based APDAs may extract different *Refined DFGs* from the same log. Also, a DFG-based APDA may discover different *Refined DFGs* from the same log depending on its hyperparameter settings (e.g., a filtering threshold). The algorithm(s) used by a DFG-based APDA to discover the *Refined DFG* from the event log and convert it into a process model may greatly affect the accuracy of an APDA. Accordingly, our framework focuses on optimizing the discovery of the *Refined DFG* rather than its conversion into a process model.

Given that a *Refined DFG* is a binary graph, it is possible to represent it in the form of a matrix as follows.

Definition 4 (DFG-Matrix) Given a *Refined DFG* $\mathcal{G} = (N, E)$ and a function $\theta : N \rightarrow [1, |N|]$,³ the DFG-Matrix is a squared matrix $X_{\mathcal{G}} \in [0, 1] \cap \mathbb{N}^{|N| \times |N|}$, where each cell $x_{i,j} = 1 \iff \exists (a_1, a_2) \in E \mid \theta(a_1) = i \wedge \theta(a_2) = j$, otherwise $x_{i,j} = 0$.

In the remaining of this paper, we refer to the *Refined DFG* as DFG for simplicity reason.

3.2 Framework overview

As shown in Fig. 1, our framework takes three inputs (in addition to the log): (i) the optimization metaheuristics; (ii) the objective function to be optimized (e.g., F-score); (iii) and the DFG-based APDA to be used for discovering a process model.

Algorithm 1 describes how our framework operates, while Fig. 2 captures the control flow representation of the Algorithm 1. First, the input event log is given to the APDA, which returns the discovered (refined) DFG and its corresponding process model (lines 1 and 2). This (refined) DFG becomes the *current DFG*, while the process model becomes the *best process model* (so far). This process model's objective function score (e.g., the F-score) is stored as the current score and the best score (lines 3 and 4). The current DFG is then given as input to the function *GenerateNeighbors*, which applies changes to the current DFG to generate a set of neighboring DFGs (line 6). The latter ones are given as input to the APDA, which returns the corresponding process models. The process models are assessed by the objective function evaluators (line 9 to 13). When the metaheuristic receives the results from the evaluators (along with the current DFG and its score), it chooses the new current DFG and updates the current score (lines 14 and 15). If the new current score is higher than the best score (line 16), it updates the best process model and the best score (lines 17 and 18). After the update, a new iteration starts, unless a termination criterion is met (e.g., a timeout, a maximum number of iterations, or a minimum threshold for the objective function). In the latter case, the framework outputs the best process model identified, i.e., the process model scoring the highest value for the objective function.

3.3 Adaptation of the optimization metaheuristics

To adapt iterative local search (ILS), tabu search (TABU), and simulated annealing (SIMA) to the problem of automated process discovery, we need to define the following

³ θ maps each node of the *Refined DFG* to a natural number.

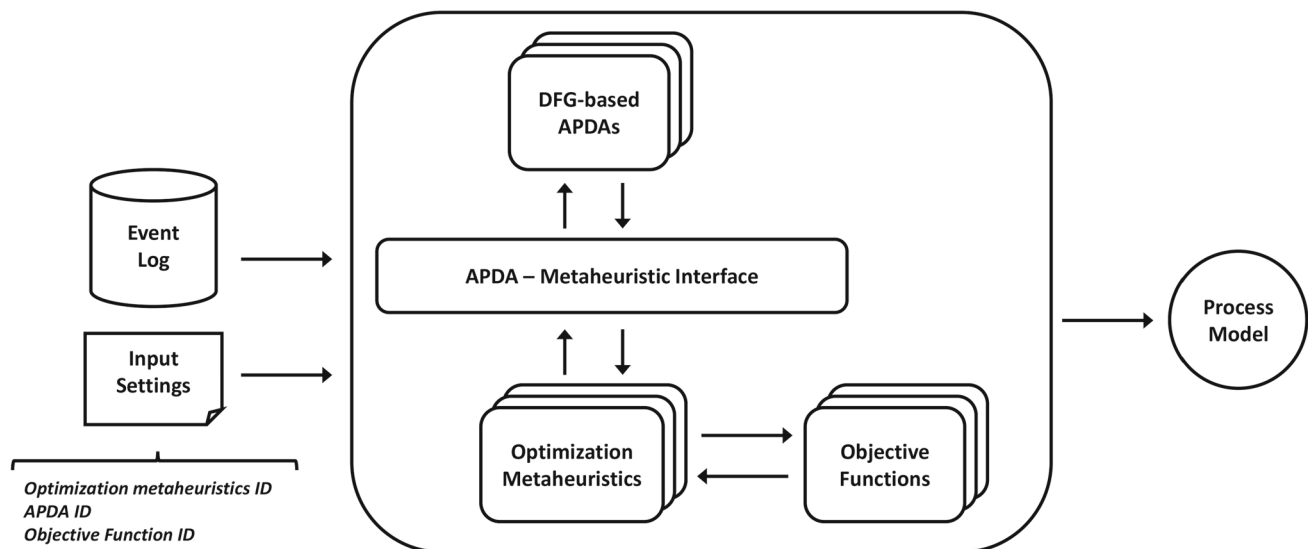


Fig. 1 Overview of our optimization framework

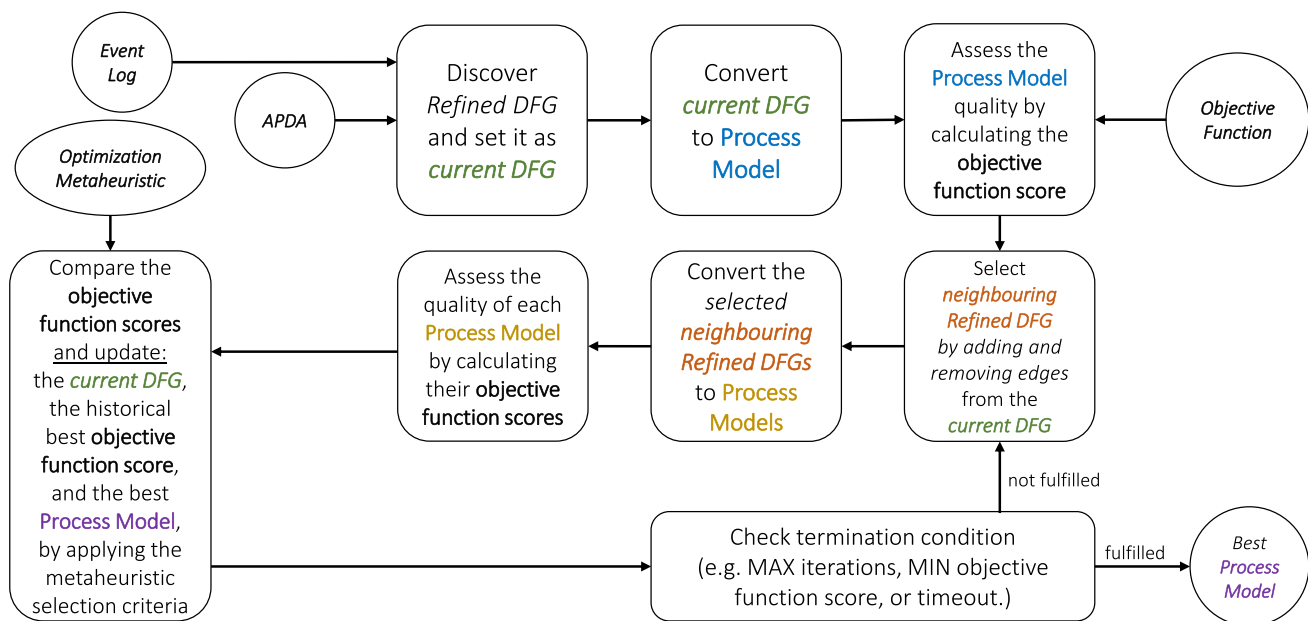


Fig. 2 Algorithm 1—control flow sketch

three concepts: (i) the problem solution space; (ii) a solution neighborhood; (iii) the objective function. These design choices influence how each of the metaheuristics navigates the solution space and escapes local minima, i.e., how to design the Algorithm 1 functions: *GenerateNeighbors* and *UpdateDFG*, resp. lines 6 and 14.

Solution space Our goal being the optimization of APDAs, we are forced to choose a solution space that fits well our context regardless the selected APDA. If we assume that the APDA is DFG-based (that is the case for the majority of the available APDAs), we can define the solution space as the set of all the DFG discoverable from the event log. Indeed, any

DFG-based APDA can generate deterministically a process model from a DFG.

Solution neighborhood Having defined the solution space as the set of all the DFG discoverable from the event log, we can refer to any element of this solution space as a DFG-Matrix. Given a DFG-Matrix, we define its neighborhood as the set of all the matrices having one different cell value (i.e., DFGs having one more/less edge). In the following, every time we refer to a DFG we assume it is represented as a DFG-Matrix.

Algorithm 1: Optimization Framework

input : Event Log \mathcal{L} , Metaheuristic ω , Objective Function \mathcal{F} , DFG-based APDA α

```

1 CurrentDFG  $\mathcal{G}_c \leftarrow \text{DiscoverRefinedDFG}(\alpha, \mathcal{L})$ ;
2 BestModel  $\hat{m} \leftarrow \text{ConvertDFGtoProcessModel}(\alpha, \mathcal{G}_c)$ ;
3 CurrentScore  $s_c \leftarrow \text{AssessQuality}(\mathcal{F}, \mathcal{L}, \hat{m})$ ;
4 BestScore  $\hat{s} \leftarrow s_c$ ;
5 while  $\text{CheckTerminationCriteria}() = \text{FALSE}$  do
6   Set  $V \leftarrow \text{GenerateNeighbors}(\mathcal{G}_c, s_c)$ ;
7   Map  $S \leftarrow \emptyset$ ;
8   Map  $M \leftarrow \emptyset$ ;
9   for  $\mathcal{G} \in V$  do
10    ProcessModel  $m \leftarrow \text{ConvertDFGtoProcessModel}(\alpha, \mathcal{G})$ ;
11    Score  $s \leftarrow \text{AssessQuality}(\mathcal{F}, \mathcal{L}, m)$ ;
12    add  $(\mathcal{G}, s)$  to  $S$ ;
13    add  $(\mathcal{G}, m)$  to  $M$ ;
14    $\mathcal{G}_c \leftarrow \text{UpdateDFG}(\omega, S, \mathcal{G}_c, s_c, \alpha, \mathcal{L})$ ;
15    $s_c \leftarrow \text{GetMapElement}(S, \mathcal{G}_c)$ ;
16   if  $\hat{s} < s_c$  then
17      $\hat{s} \leftarrow s_c$ ;
18      $\hat{m} \leftarrow \text{GetMapElement}(M, \mathcal{G}_c)$ ;
19 return  $\hat{m}$ ;

```

Objective function It is possible to define the objective function as any function assessing one of the four quality dimensions for discovered process models (introduced in Sect. 2). However, being interested in optimizing the APDAs to discover the most accurate process model, in our optimization framework instantiations we refer to the objective function as the F-score of fitness and precision. Furthermore, we remark that our framework could operate also with objective functions that take into account multiple quality dimensions striving for a trade-off, e.g., F-score and model complexity, provided the multiple quality dimensions can be combined into a unique objective function.

Having defined the solution space, a solution neighborhood, and the objective function, we can turn our attention on

how ILS, TABU, and SIMA navigate the solution space. ILS, TABU, and SIMA share similar traits in solving an optimization problem, especially when it comes to the navigation of the solution space. Given a problem and its solution space, any of these three S-metaheuristics starts from a (random) solution, discovers one or more neighboring solutions, and assesses them with the objective function to find a solution that is better than the current one. If a better solution is found, it is chosen as the new current solution and the metaheuristic performs a new neighborhood exploration. If a better solution is not found, e.g., the current solution is locally optimal, the three metaheuristics follow different approaches to escape the local optimum and continue the solution space exploration. Algorithm 1 orchestrates and facilitates the parts of this procedure shared by the three metaheuristics. However, we must define the functions *GenerateNeighbors* (GNF) and *UpdateDFG* (UDF).

The GNF receives in input a solution of the solution space, i.e., a DFG, and it generates a set of neighboring DFGs. By definition, GNF is independent from the metaheuristic and it can be as simple or as elaborate as we demand. An example of a simple GNF is a function that randomly selects neighboring DFGs turning one cell of the input DFG-Matrix to 0 or to 1. While, an example of an elaborate GNF is a function that accurately selects neighboring DFGs relying on the feedback received from the objective function assessing the input DFG, as we show in Sect. 3.4.

The UDF (captured in Algorithm 2) is the core of our optimization framework, and it implements the metaheuristic itself. The UDF receives in input the selected metaheuristic (ω), the neighboring DFGs and their corresponding objective function scores (S), the current DFG (\mathcal{G}_c), the current score

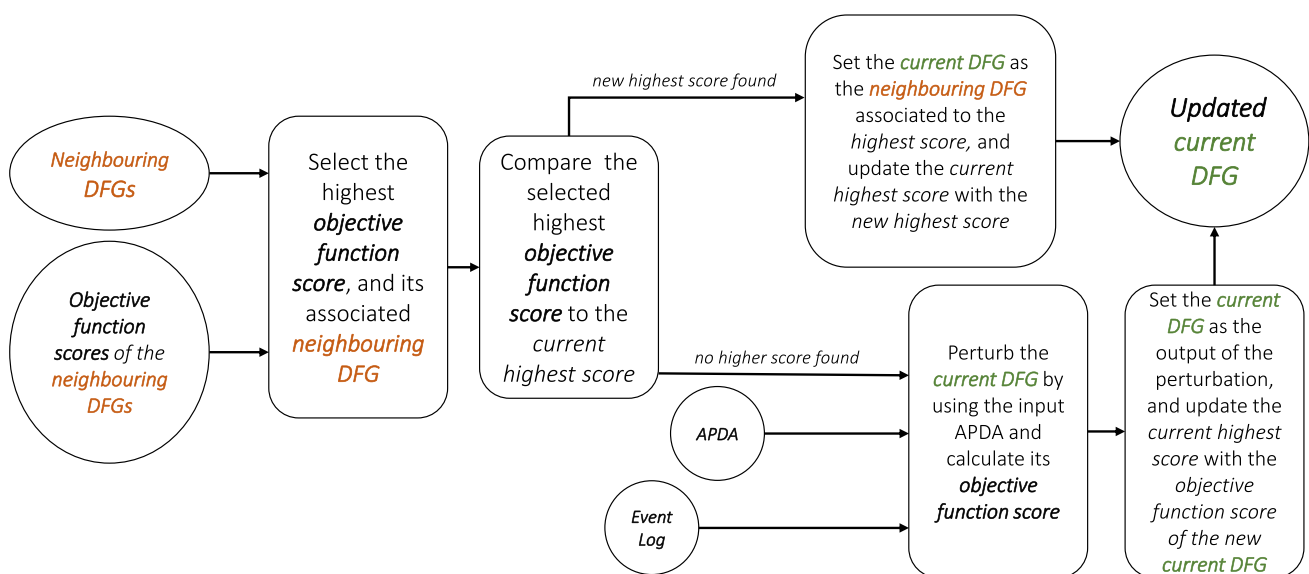


Fig. 3 UDF when selecting ILS as optimization metaheuristic

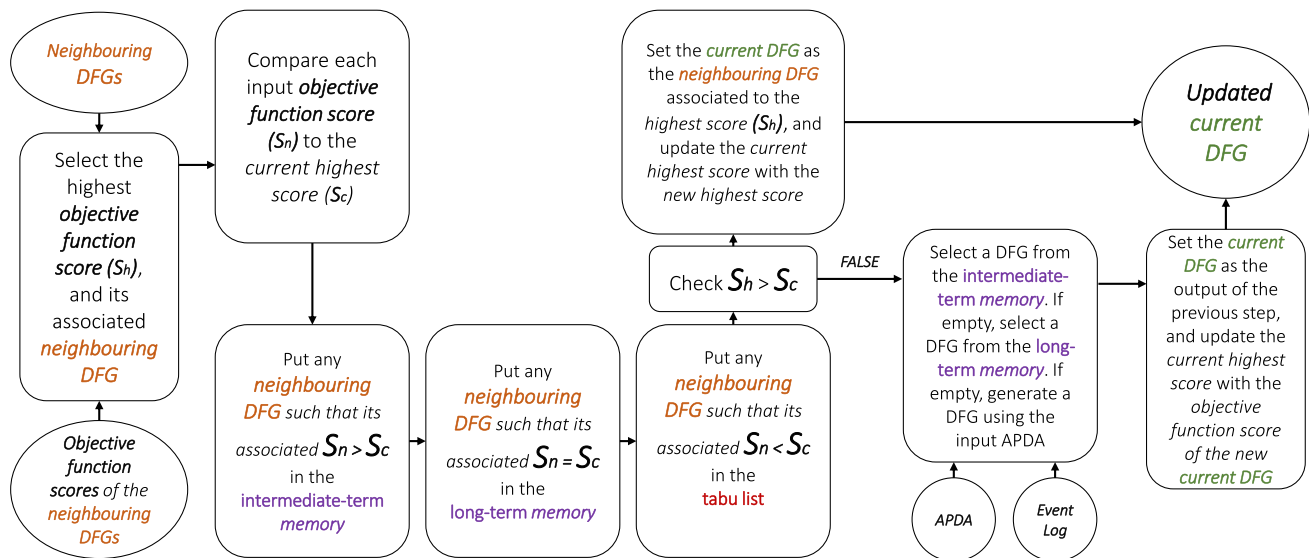


Fig. 4 UDF when selecting TABU as optimization metaheuristic

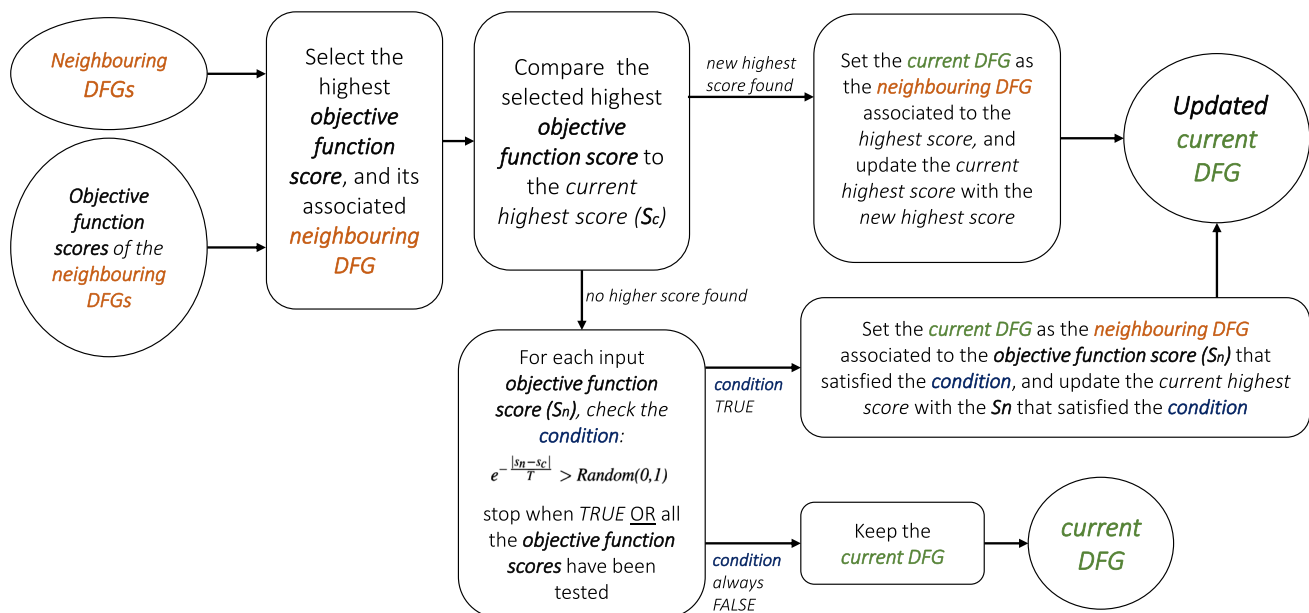


Fig. 5 UDF when selecting SIMA as optimization metaheuristic

(s_c), the APDA (α), and the event log (\mathcal{L}). Then, we can differentiate two cases: (i) among the input neighboring DFGs there is at least one having a higher objective function score than the current; (ii) none of the input neighboring DFGs has a higher objective function score than the current. In the first case, UDF always outputs the DFG having the highest score regardless of the selected metaheuristic (see Algorithm 2,

lines 4, 11, and 33—respectively, for ILS, TABU, and SIMA). In the second case, the current DFG may be a local optimum, and each metaheuristic escapes it with a different strategy. Figures 3, 4, and 5 show the high-level control flow of how ILS, TABU, and SIMA update the current DFG (that is, the UDF—Algorithm 2).

Algorithm 2: Update DFG Function (UDF)

```

input : Metaheuristic  $\omega$ , Neighboring-DFG Scores Map  $S$ , Current DFG  $\mathcal{G}_c$ , Current best score  $s_c$ , DFG-based APDA  $\alpha$ , Event Log  $\mathcal{L}$ 

1  $\mathcal{G}_b \leftarrow \mathcal{G}_c$ ;
2 if  $\omega = \text{ILS}$  then
3   for  $\mathcal{G} \in \text{GetMapKeys}(S)$  do
4     if  $\text{GetMapElement}(S, \mathcal{G}) > s_c$  then
5        $\mathcal{G}_b \leftarrow \mathcal{G}$ ;
6        $s_c \leftarrow \text{GetMapElement}(S, \mathcal{G})$ ;
7   if  $\mathcal{G}_b = \mathcal{G}_c$  then  $\mathcal{G}_b \leftarrow \text{PerturbDFG}(\alpha, \mathcal{G}_c)$ ;
8 if  $\omega = \text{TABU}$  then
9    $s_{tmp} \leftarrow s_c$ ;
10  for  $\mathcal{G} \in \text{GetMapKeys}(S)$  do
11    if  $\text{GetMapElement}(S, \mathcal{G}) > s_c$  then
12       $\mathcal{G}_b \leftarrow \mathcal{G}$ ;
13       $s_c \leftarrow \text{GetMapElement}(S, \mathcal{G})$ ;
14  remove  $\mathcal{G}_b$  from  $S$ ;
15   $s_c \leftarrow s_{tmp}$ ;
16  for  $\mathcal{G} \in \text{GetMapKeys}(S)$  do
17    if  $\text{GetMapElement}(S, \mathcal{G}) \geq s_c$  then
18      if  $\text{GetMapElement}(S, \mathcal{G}) > s_c$  then
19        add  $\mathcal{G}$  to  $\text{IntermediateTermMemory}(\omega)$ ;
20      else
21        add  $\mathcal{G}$  to  $\text{LongTermMemory}(\omega)$ ;
22    else
23      add  $\mathcal{G}$  to  $\text{TabuList}(\omega)$ ;
24  add  $\mathcal{G}_c$  to  $\text{TabuList}(\omega)$ ;
25  if  $\mathcal{G}_b = \mathcal{G}_c$  then
26    if  $\text{IntermediateTermMemory}(\omega) \neq \emptyset$  then
27       $\mathcal{G}_b \leftarrow \text{GetLastIntermediateTermMemoryElement}(\omega)$ ;
28    else
29      if  $\text{LongTermMemory}(\omega) \neq \emptyset$  then  $\mathcal{G}_b \leftarrow$ 
30         $\text{GetLastLongTermMemoryElement}(\omega)$ ;
31       $\mathcal{G}_b \leftarrow \text{DiscoverDFG}(\alpha, \mathcal{L})$ ;
31 if  $\omega = \text{SIMA}$  then
32   for  $\mathcal{G} \in \text{GetMapKeys}(S)$  do
33     if  $\text{GetMapElement}(S, \mathcal{G}) > s_c$  then
34        $\mathcal{G}_b \leftarrow \mathcal{G}$ ;
35        $s_c \leftarrow \text{GetMapElement}(S, \mathcal{G})$ ;
36   if  $\mathcal{G}_b = \mathcal{G}_c$  then
37     for  $\mathcal{G} \in \text{GetMapKeys}(S)$  do
38        $s_n \leftarrow \text{GetMapElement}(S, \mathcal{G})$ ;
39       if  $e^{-\frac{|s_n - s_c|}{T}} > \text{Random}(0, 1)$  then
40          $\mathcal{G}_b \leftarrow \mathcal{G}$ ;
41 return  $\mathcal{G}_b$ ;

```

Iterative Local Search applies the simplest strategy, it perturbs the current DFG (Algorithm 2, line 7). The perturbation is meant to alter the DFG in such a way to escape the local optimum, e.g., randomly adding and removing multiple edges from the current DFG. The perturbed DFG is the output of the UDF.

Tabu Search relies on its three memories to escape a local optimum (Algorithm 2, line 25 to 30). The short-term memory (a.k.a. Tabu-list), which contains DFG that must not be explored further. The intermediate-term memory, which contains DFGs that should lead to better results and, therefore, should be explored in the near future. The long-term memory, which contains DFGs (with characteristics) that have been seen multiple times and, therefore, not to explore in the near

future. TABU updates the three memories each time the UDF is executed. Given the set of neighboring DFGs and their respective objective function scores (see Algorithm 1, map S), TABU adds each DFG to a different memory. DFGs worsening the objective function score are added to the Tabu-list. DFGs improving the objective function score, yet less than another neighboring DFG, are added to the intermediate-term memory. DFGs that do not improve the objective function score are added to the long-term memory. Also, the current DFG is added to the Tabu-list, it being already explored. When TABU does not find a better DFG in the neighborhood of the current DFG, it returns the latest DFG added to the intermediate-term memory. If the intermediate-term memory is empty, TABU returns the latest DFG added to the long-term memory. If both these memories are empty, TABU requires a new (random) DFG from the APDA, and outputs its DFG.

Simulated Annealing avoids getting stuck in a local optimum by allowing the selection of DFGs worsening the objective function score (Algorithm 2, line 36 to 40). In doing so, SIMA explores areas of the solution space that other S-metaheuristics do not. When a better DFG is not found in the neighborhood of the current DFG, SIMA analyzes one neighboring DFG at a time. If this neighbor does not worsen the objective function score, SIMA outputs it. Instead, if the neighboring DFG worsens the objective function score, SIMA outputs it with a probability of $e^{-\frac{|s_n - s_c|}{T}}$, where s_n and s_c are the objective function scores of, (respectively), the neighboring DFG and the current DFG, and the temperature T is an integer that converges to zero as a linear function of the maximum number of iterations. The temperature is fundamental to avoid updating the current DFG with a worse one if there would be no time to recover from the worsening (i.e., too few iterations left for continuing the exploration of the solution space from the worse DFG).

3.4 Framework instantiation

To assess our framework, we instantiated it for three APDAs: Split Miner [9], Fodina [34], and Inductive Miner [24]. These three APDAs are all DFG-based, and they are representatives of the state of the art. In fact, the latest APDAs literature review and benchmark [8] showed that Fodina, Split Miner, and Inductive Miner outperformed other APDAs when their hyperparameters were optimized via a brute-force approach. Therefore, we decided to focus on those DFG-based APDAs that would benefit the most from the application of our optimization framework.

To complete the instantiation of our framework for any concrete DFG-based APDA, it is necessary to implement an interface that allows the metaheuristics to interact with the APDA (as discussed above). Such an interface should

provide four functions: *DiscoverDFG* and *ConvertDFGtoProcessModel* (see Algorithm 1), the *Restart Function* (RF) for TABU, and the *Perturbation Function* (PF) for ILS.

The first two functions, *DiscoverDFG* and *ConvertDFGtoProcessModel*, are inherited from the DFG-based APDA, in our case Split Miner, Fodina, and Inductive Miner. We note that Split Miner and Fodina receive as input parameter settings that can vary the output of the *DiscoverDFG* function. Precisely, *Split Miner* has two parameters: the noise filtering threshold, used to drop infrequent edges in the DFG, and the parallelism threshold, used to determine which potential parallel relations between activities are used when discovering the process model from the DFG. While, *Fodina* has three parameters: the noise filtering threshold, similar to the one of Split Miner, and two threshold to detect, respectively, self-loops and short-loops in the DFG. Instead, the DFG-based variant of Inductive Miner [24] that we integrated in our optimization framework does not receive any input parameters.

To discover the initial DFG (Algorithm 1, line 1) with Split Miner, default parameters are used.⁴ We removed the randomness for discovering the initial DFG because most of the times, the DFG discovered by Split Miner with default parameters is already a good solution [9], and starting the solution space exploration from this latter can reduce the total exploration time.

Similarly, if Fodina is the selected APDA, the initial DFG (Algorithm 1, line 1) is discovered using the default parameters of Fodina,⁵ even though there is no guarantee that the default parameters allow Fodina to discover a good starting solution [8]. Yet, this design choice is less risky than randomly choosing the values of the input parameters in order to discover the initial DFG, because it is likely Fodina would discover unsound models when randomly tuned, given that it does not guarantee soundness.

On the other hand, Inductive Miner [24] does not apply any manipulation to the discovered initial DFG. In this case, we pseudorandomly generate an initial DFG starting from a given seed, to ensure determinism. Differently than the case of Fodina, this is a suitable design choice for Inductive Miner, because it always guarantees block-structured sound process models, regardless of the DFG.

Function RF is very similar to *DiscoverDFG*, since it requires the APDA to output a DFG. The only difference is that RF must output a different DFG every time it is executed. We adapted the *DiscoverDFG* function of Split Miner and Fodina to output the DFG discovered with default parameters the first time it is executed, and a DFG discovered with pseudorandom parameters for the following executions.

The case of Inductive Miner is simpler, because the *DiscoverDFG* function always returns a pseudorandom DFG. Consequently, we mapped RF to the *DiscoverDFG* function.

Finally, function PF can be provided either by the APDA (through the interface) or by the metaheuristic. However, PF can be more effective when not generalised by the metaheuristic, allowing the APDA to apply different perturbations to the DFGs, taking into account how the APDA converts the DFG to a process model. We chose a different PF for each of the three APDAs.

- *Split Miner PF* We invoke Split Miner's concurrency oracle to extract the possible parallelism relations in the log using a randomly chosen parallelism threshold. For each new parallel relation discovered that is not present in the current solution, two edges are removed from the DFG, whilst, for each deprecated parallel relation, two edges are added to the DFG.
- *Fodina PF* Given the current DFG, we analyze its self-loops and short-loops relations using random loop thresholds. As a result, a new DFG is generated where a different set of edges is retained as self-loops and short-loops.
- *Inductive Miner PF* Since Inductive Miner does not perform any manipulation on the DFG, we could not determine an efficient way to perturb the DFG. Thus, we set PF = RF, so that instead of perturbing the current DFG, a new random DFG is generated. This variant of the ILS is called Repeated local search (RLS). In the evaluation reported in Sect. 4, we use only RLS for Inductive Miner, and both ILS and RLS for Fodina and Split Miner.

To complete the instantiation of our framework, we need to set an objective function. With the goal of optimizing the accuracy of the APDAs, we chose as objective function the F-score of fitness and precision. Among the existing measures of fitness and precision, we selected the Markovian fitness and precision presented in [5,6].⁶ The rationale for this choice is that these measures of fitness and precision are the fastest to compute among state-of-the-art measures [5,6]. Furthermore, these measures indicate what edges could be added to or removed from the DFG to improve the fitness or precision of the model. This feedback allows us to design an effective GNF.

In the instantiation of our framework, the objective function's output is a data structure composed of: the Markovian fitness and precision of the model, the F-score, and the mismatches between the model and the event log identified during the computation of the Markovian fitness and precision, i.e., the sets of edges that could be added to improve

⁴ The default parameters of Split Miner are the most likely to yield the best results [9].

⁵ The default parameters of Fodina are the most likely to yield the best results [34].

⁶ We used the Boolean function variant with order $k = 5$.

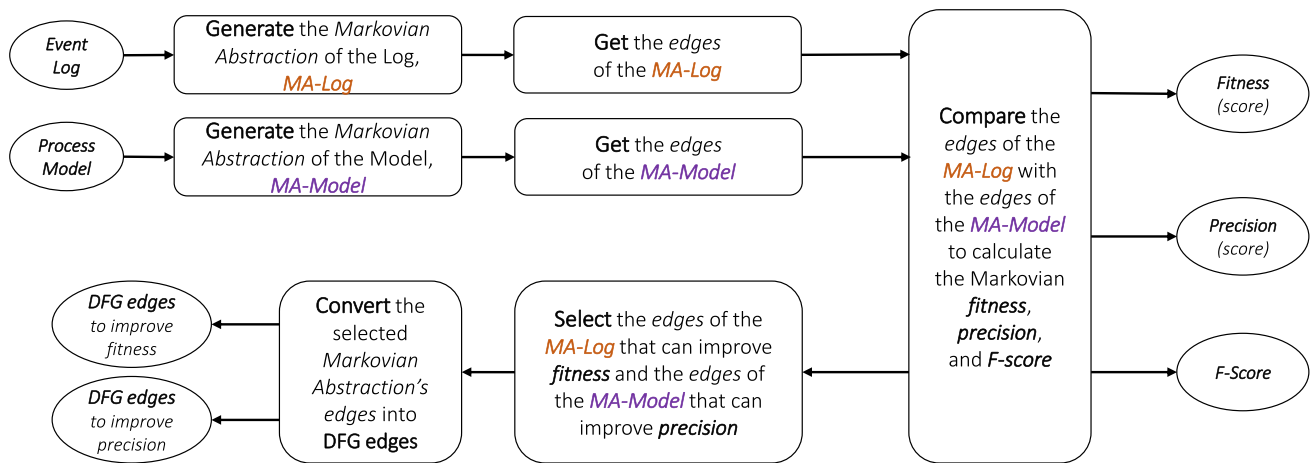


Fig. 6 Algorithm 3—control flow sketch

fitness or removed to improve precision. Algorithm 3 illustrates how we build this data structure, its high-level control flow sketch is captured in Fig. 6.

Given an event log and a process model, we generate their respective Markovian abstractions by applying the method described in [5] (lines 1 and 2). We recall that the Markovian abstraction of the log/model is a graph, where each edge represents a subtrace⁷ observed in the log/model. Next, we collect all the edges of the Markovian abstraction of the log and of the model into two different sets: E_l and E_m (lines 3 and 4). These two sets are used to determine the Markovian fitness and precision of the process model [5], by applying the formula in lines 4 and 10. We note that the edges in E_l that cannot be found in E_m (set E_{df} , line 6) represent subtraces of the log that cannot be found in the process model. Vice-versa, the edges in E_m that cannot be found in E_l (set E_{dp} , line 11) represent subtraces of the process model that cannot be found in the log. We analyze these subtraces to detect directly-follows relations, i.e., DFG edges (lines 9 and 14), that can be added or removed from the DFG that generated the process model in order to either improve fitness or precision. Precisely, the DFG edges that can be added to improve fitness are those captured by the directly-follows relations that we can find in the Markovian abstraction edges in set E_{df} . On the other hand, the edges that can be removed to improve precision are those captured by the directly-follows relations that we can find in the Markovian abstraction edges in set E_{dp} . Once these edges to be added or removed are identified (sets E_f and E_p), we can output the final data structure, which comprises the Markovian fitness and precision, their F-score, and the two sets E_f and E_p .

⁷ The length of the subtrace is determined by the order of the Markovian abstraction, in our case $k = 5$.

Algorithm 3: Assess Quality with Markovian Accuracy

```

input : Event Log  $\mathcal{L}$ , Model  $m$ 
1 MarkovianAbstraction  $L \leftarrow \text{GenerateMarkovianAbstraction}(\mathcal{L})$ ;
2 MarkovianAbstraction  $P \leftarrow \text{GenerateMarkovianAbstraction}(m)$ ;
3 Set  $E_l \leftarrow \text{GetEdges}(L)$ ;
4 Set  $E_m \leftarrow \text{GetEdges}(P)$ ;
5  $f \leftarrow 1 - \frac{|E_l \setminus E_m|}{|E_l|}$ ;
6 Set  $E_{df} \leftarrow E_l \setminus E_m$ ;
7 Set  $E_f \leftarrow \emptyset$ ;
8 for  $e \in E_{df}$  do
9    $E_f \leftarrow E_f \cup \text{ExtractDFGEdges}(e)$ ;
10  $p \leftarrow 1 - \frac{|E_m \setminus E_l|}{|E_m|}$ ;
11  $E_{dp} \leftarrow E_m \setminus E_l$ ;
12 Set  $E_p \leftarrow \emptyset$ ;
13 for  $e \in E_{dp}$  do
14    $E_p \leftarrow E_p \cup \text{ExtractDFGEdges}(e)$ ;
15  $s \leftarrow \frac{2 \cdot f \cdot p}{f + p}$ ;
16 return  $(s, f, p, E_f, E_p)$ ;
  
```

Given the above objective function's output, our GNF is described in Algorithm 4, while Fig. 7 captures its high-level control flow sketch.

Algorithm 4: Generate Neighbors Function (GNF)

```

input : CurrentDFG  $\mathcal{G}_c$ , CurrentMarkovianScore  $s_c$ , Integer  $size_n$ 
1 if  $\text{getFitnessScore}(s_c) > \text{getPrecisionScore}(s_c)$  then
2   Set  $E_m \leftarrow \text{getEdgesForImprovingPrecision}(s_c)$ ;
3 else
4   Set  $E_m \leftarrow \text{getEdgesForImprovingFitness}(s_c)$ ;
5 Set  $N \leftarrow \emptyset$ ;
6 while  $E_m \neq \emptyset \wedge |N| \neq size_n$  do
7   Edge  $e \leftarrow \text{getRandomElement}(E_m)$ ;
8   NeighboringDFG  $\mathcal{G}_n \leftarrow \text{copyDFG}(\mathcal{G}_c)$ ;
9   if  $\text{getFitnessScore}(s_c) > \text{getPrecisionScore}(s_c)$  then
10    if  $\text{canRemoveEdge}(\mathcal{G}_n, e)$  then add  $\mathcal{G}_n$  to  $N$ ;
11  else
12    addEdge( $\mathcal{G}_n, e$ );
13    add  $\mathcal{G}_n$  to  $N$ ;
14 return  $N$ ;
  
```

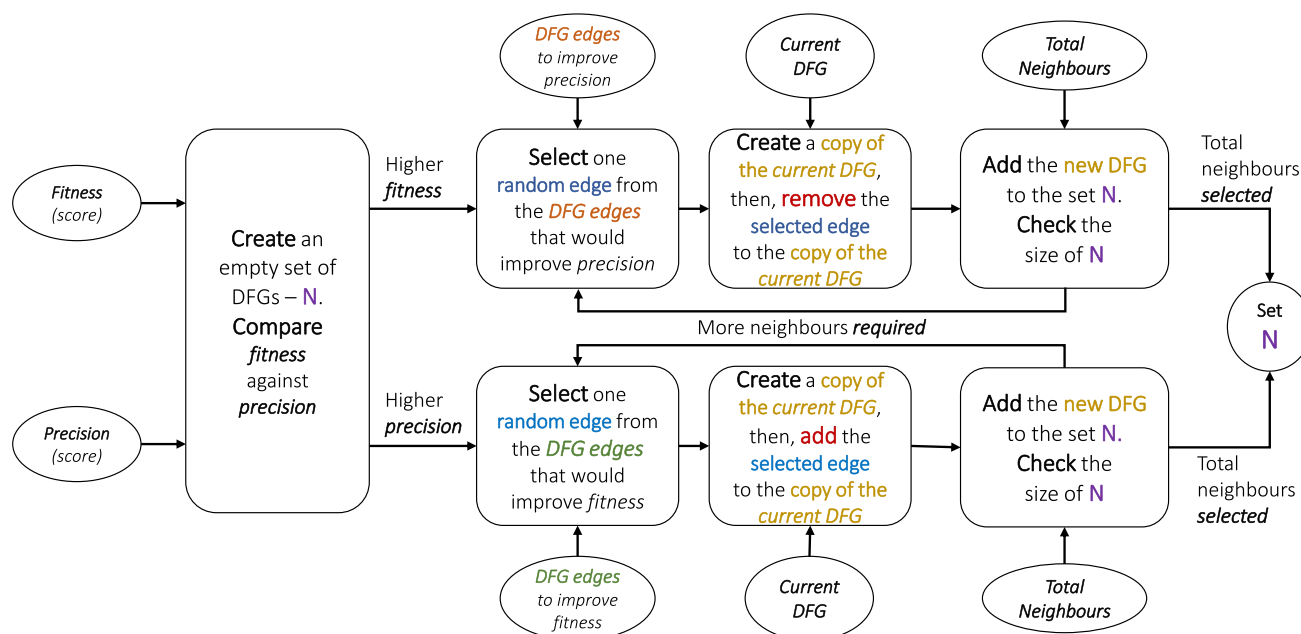


Fig. 7 Algorithm 4—control flow sketch

This function receives as input the current DFG (\mathcal{G}_c), its objective function score (the data structure s_c), and the number of neighbors to generate ($size_n$). If fitness is greater than precision, we retrieve from s_c the set of edges (E_m) that could be removed from \mathcal{G}_c to improve its precision (line 2). Conversely, if precision is greater than fitness, we retrieve from s_c the set of edges (E_m) that could be added to \mathcal{G}_c to improve its fitness (line 4). The reasoning behind this design choice is that, given that our objective function is the F-score, it is preferable to increase the lowest of the two measures (precision or fitness). That is, if the fitness is lower, we increase fitness, and conversely if the precision is lower we increase precision. Once we have E_m , we randomly select one edge from it, generate a copy of the current DFG (\mathcal{G}_n), and either remove or add the randomly selected edge according to the accuracy measure we want to improve (precision or fitness), see lines 7 to 13. If the removal of an edge generates a disconnected \mathcal{G}_n , we do not add this latter to the neighbors set (N), line 10. We keep iterating over E_m until the set is empty (i.e., no mismatching edges are left) or N reaches its maximum size (i.e., $size_n$). We then return N . The algorithm ends when the maximum execution time or the maximum number of iterations is reached.

4 Evaluation

We implemented the proposed optimization framework as a Java command-line application.⁸ This tool uses Split Miner,

⁸ Available under the label “Optimization Framework for Automated Process Discovery” at <http://apromore.org/platform/tools>.

Fodina, and Inductive Miner as the underlying APDAs, and the Markovian accuracy F-score as the objective function (cf. Sect. 3.4). Using this implementation, we undertook to empirically evaluate the magnitude of improvements in accuracy delivered by different instantiations of the framework.

4.1 Dataset, quality measures, and experimental setup

For our evaluation, we used the dataset of the benchmark of automated process discovery approaches in [8], which to the best of our knowledge is the most recent benchmark on this topic. This dataset includes twelve public logs and eight private logs. The public logs originate from the 4TU Centre for Research Data and include the *BPI Challenge* (BPIC) logs (2012–17), the *Road Traffic Fines Management Process* (RTFMP) log and the *SEPSIS* log. These logs record executions of business processes from a variety of domains, e.g., healthcare, finance, government, and IT service management. The eight proprietary logs are sourced from several companies in the education, insurance, IT service management, and IP management domains.

Table 1 reports the characteristics of the logs. The dataset comprises simple logs (e.g., BPIC13_{cp}) and very complex ones (e.g., SEPSIS, PRT2) in terms of percentage of distinct traces, and both small logs (e.g., BPIC13_{cp} and SEPSIS) and large ones (e.g., BPIC17 and PRT9) in terms of total number of events.

From each of these logs, we discovered 16 process models by applying the following techniques:

Table 1 Descriptive statistics of the real-life logs (public and proprietary)

Log	BPIC12	BPIC13 _{cp}	BPIC13 _{inc}	BPIC14 _f	BPIC15 _{if}	BPIC15 _{2f}	BPIC15 _{3f}	BPIC15 _{4f}	BPIC15 _{5f}
Total traces	13,087	1487	7554	41,353	902	681	1369	860	975
Dist. traces (%)	33.4	12.3	20	36.1	32.7	61.7	60.3	52.4	45.7
Total events	262,200	6660	65,533	369,485	21,656	24,678	43,786	29,403	30,030
Dist. events	36	7	13	9	70	82	62	65	74
Tr. length (min)	3	1	1	3	5	4	4	5	4
(avg)	20	4	9	9	24	36	32	34	31
(max)	175	35	123	167	50	63	54	54	61

Log	BPIC17 _f	RTFMP	SEPSIS	PRT1	PRT2	PRT3	PRT4	PRT6	PRT7	PRT9	PRT10
Total traces	21,861	150,370	1050	12,720	1182	1600	20,000	744	2000	787,657	43,514
Dist. traces (%)	40.1	0.2	80.6	8.1	97.5	19.9	29.7	22.4	6.4	0.01	0.01
Total events	714,198	561,470	15,214	75,353	46,282	13,720	166,282	6011	16,353	1,808,706	78,864
Dist. events	41	11	16	9	9	15	11	9	13	8	19
Tr. length (min)	11	2	3	2	12	6	6	7	8	1	1
(avg)	33	4	14	5	39	8	8	8	8	2	1
(max)	113	2	185	64	276	9	36	21	11	58	15

- Split Miner with default parameters (SM);
- Split Miner with hyper-parameter optimization⁹ (HPO_{sm});
- Split Miner optimized with our framework using the following optimization metaheuristics: RLS_{sm}, ILS_{sm}, TABU_{sm}, SIMA_{sm};
- Fodina with default parameters (FO);
- Fodina with hyper-parameter optimization⁹ (HPO_{fo});
- Fodina optimized with our framework using the following optimization metaheuristics: RLS_{fo}, ILS_{fo}, TABU_{fo}, SIMA_{fo};
- Inductive Miner IM_d;
- Inductive Miner optimized with our framework using the following optimization metaheuristics: RLS_{imd}, TABU_{imd}, SIMA_{imd}.⁹

For each of the above metaheuristics, we set the maximum execution time to five minutes and the maximum number of iterations to 50. The same timeout was also applied to the hyper-parameter optimizations.

For each of the discovered models, we measured fitness, precision, complexity, and execution time. For measuring fitness and precision, we adopted two different sets of measures. The first set of measures is based on *alignments*, computing fitness, and precision with the approaches proposed by Adriansyah et al. [1,2] (alignment-based accuracy). *Alignment-based fitness* selects for each trace in the log, the closest trace recognized by the process model, and measures the minimal number of error-corrections required to align these two traces (a.k.a. *minimal alignment cost*). The final fitness score is equal to one minus the normalized sum of the minimal alignment cost between each trace in the log and the closest corresponding trace recognized by the model. *Alignment-based precision* builds a *prefix automaton* from the event log, then it replays the process model behavior on top of the log prefix automation (with the aid of alignments) and counts the number of times that the model can perform a move that the prefix automaton cannot. Each of these mismatching moves is called an *escaping edge*. The final value of precision is function of the number of detected escaping edges. For more details regarding the alignment-based fitness and precision, we refer to the corresponding studies [1,2].

The second set of measures is based on *Markovian abstractions*, computing fitness, and precision with the approaches in [5]. The *Markovian fitness* generates a Markovian abstraction from the behavior recorded in the event log and a Markovian abstraction from the behavior allowed by the process model. As mentioned in the previous section, a Markovian abstraction is a graph where each node represents a subtrace of a fixed length. The *Markovian fitness* relies on a graph comparison algorithm [23] to identify the edges of the Markovian abstraction generated from the log that do not

⁹ Using the Markovian accuracy F-score as objective function.

appear in the Markovian abstraction generated from the process model. Similarly, the *Markovian precision* is calculated by identifying (via the same graph comparison algorithm [23]) the edges of the Markovian abstraction of the process model that do not appear in the Markovian abstraction of the log. For more details regarding the alignment-based fitness and precision, we refer to the corresponding study [5].

For assessing the complexity of the models we relied on *size*, control-flow complexity (CFC), and structuredness. *Size* is the total number of nodes of a process model; *Control flow complexity* (CFC) is the amount of branching induced by the split gateways in a process model; *Structuredness* is the percentage of nodes located inside a single-entry single-exit fragment of a process model.

Note that we did not measure the generalization of the discovered process models because available generalization measures assess the capability of an APDA to generalise the behavior recorded in the event log during the discovery of a process model, and they do not assess the generalization of the process model itself [32]. However, this should not be seen as a limitation of this study, since our objective is to analyze the benefits yielded by our optimization framework in terms of F-score (through fitness and precision).

We used the results of these measurements to compare the quality of the models discovered by each baseline APDA (SM, FO, IM_d) against the quality of the models discovered by the respective optimized approaches.

All the experiments were performed on an Intel Core i5-6200U@2.30GHz with 16GB RAM running Windows 10 Pro (64-bit) and JVM 8 with 14GB RAM (10GB Stack and 4GB Heap). The framework implementation, the batch tests, the results, and all the (public) models discovered during the experiments are available for reproducibility purposes at <https://doi.org/10.6084/m9.figshare.11413794>.

4.2 Split Miner

Tables 2 and 3 show the results of our comparative evaluation for Split Miner. Each row reports the quality of each discovered process model in terms of accuracy (both alignment-based and Markovian), complexity, and discovery time. We held out from the tables four logs: BPIC13_{cp}, BPIC13_{inc}, BPIC17, and PRT9. For these logs, none of the metaheuristics could improve the accuracy of the model already discovered by SM. This is due to the high fitness score achieved by SM in these logs. By design, our metaheuristics try to improve precision by removing edges, but in these four cases, no edge could be removed without compromising the structure of the model (i.e., the model would become disconnected).

For the remaining 16 logs, all the metaheuristics consistently improved the Markovian F-score over that achieved by SM. Also, all the metaheuristics performed better than

HPO_{sm}, except in two cases (BPIC12 and PRT1). Overall, the most effective optimization metaheuristic was ILS, which delivered the highest Markovian F-score nine times out of 16, followed by SIMA_{sm} (eight times), RLS_{sm} and TABU_{sm} (six times each). We note, however, that the F-score difference between the four metaheuristics is small (in the order of one to two percentage points).

Despite the fact that the objective function of the metaheuristics was the Markovian F-score, all four metaheuristics also optimized in half of the cases the alignment-based F-score. This is due to the fact that any improvement in the Markovian fitness translates into an improvement in the alignment-based fitness. This does not hold for precision. The result highlights the partial correlation between alignment-based and Markovian measures, already discussed in the previous section.

By close inspection to the complexity of the models, we note that most of the times (nine cases out of 16) the F-score improvement achieved by the metaheuristics comes at the cost of size and CFC. This is expected, since SM tends to discover models with higher precision than fitness [9]. To improve the F-score, new behavior is added to the model in the form of new edges (note that new nodes are never added); this leads to new gateways and consequently to higher size and CFC. On the other hand, when precision is lower than fitness, and thus the metaheuristic aims to increase the value of precision to improve the overall F-score, the result is the opposite: the model complexity reduces as edges are removed. This is the case of the RTFMP and PRT10 logs. Supporting examples of these two possible scenarios are Figs. 8 and 9. Figure 8 shows the models discovered by SIMA_{sm} and SM from the BPIC14_f log, where the model discovered by SIMA_{sm} is more complex than that obtained with SM because it was necessary to improve its fitness (adding edges). While Fig. 9 shows the models discovered by SIMA_{sm} and SM from the RTFMP log, where the model discovered by SIMA_{sm} is simpler than that obtained with SM because it was necessary to improve the precision (removing edges).

Comparing the results obtained by the metaheuristics with HPO_{sm}, we can see that our approach allows us to discover models that cannot be discovered simply by tuning the parameters of SM. This relates to the solution space exploration. Indeed, HPO_{sm} can only explore a limited number of solutions (DFGs), i.e., those that can be generated by underlying APDA (SM in this case) by varying its parameters. In contrast, the metaheuristics go beyond the solution space of HPO_{sm} by exploring new DFGs in a pseudorandom manner.

In terms of execution times, the four metaheuristics perform similarly, having an average discovery time close to 150s. While this is considerably higher than the execution

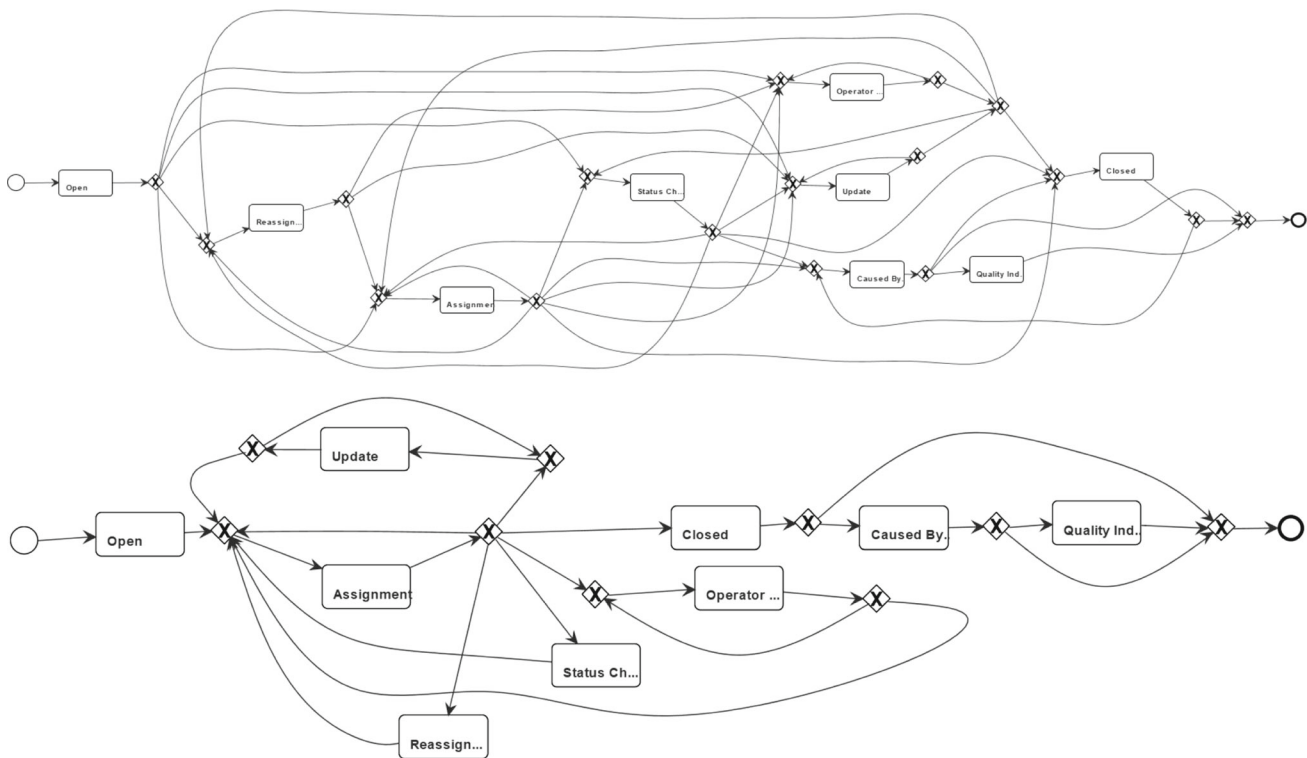
Table 2 Comparative evaluation results for the public logs—Split Miner

Event log	Discovery approach	Align. Acc.			Markov. Acc.			Complexity			Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.	
BPIC12	SM	0.963	0.520	0.675	0.818	0.139	0.238	51	41	0.69	3.2
	HPO _{sm}	0.781	0.796	0.788	0.575	0.277	0.374	40	17	0.58	4295.8
	RLS _{sm}	0.921	0.671	0.776	0.586	0.247	0.348	49	31	0.90	159.3
	ILS _{sm}	0.921	0.671	0.776	0.586	0.247	0.348	49	31	0.90	159.4
	TABU _{sm}	0.921	0.671	0.776	0.586	0.247	0.348	49	31	0.90	140.7
	SIMA _{sm}	0.921	0.671	0.776	0.586	0.247	0.348	49	31	0.90	151.1
BPIC14 _f	SM	0.772	0.881	0.823	0.150	1.000	0.262	20	14	1.00	0.8
	HPO _{sm}	0.852	0.857	0.855	0.449	1.000	0.619	22	16	0.59	575.8
	RLS _{sm}	1.000	0.771	0.871	1.000	0.985	0.992	28	34	0.54	139.0
	ILS _{sm}	1.000	0.771	0.871	1.000	0.985	0.992	28	34	0.54	151.3
	TABU _{sm}	0.955	0.775	0.855	0.856	0.999	0.922	26	31	0.69	154.7
	SIMA _{sm}	1.000	0.771	0.871	1.000	0.985	0.992	28	34	0.54	140.3
BPIC15 _{1f}	SM	0.899	0.871	0.885	0.701	0.726	0.713	111	45	0.51	0.7
	HPO _{sm}	0.962	0.833	0.893	0.804	0.670	0.731	117	55	0.45	1242.3
	RLS _{sm}	0.925	0.839	0.880	0.774	0.803	0.788	124	63	0.39	163.6
	ILS _{sm}	0.925	0.839	0.880	0.774	0.803	0.788	124	63	0.39	166.8
	TABU _{sm}	0.948	0.843	0.892	0.774	0.805	0.789	125	64	0.33	187.2
	SIMA _{sm}	0.920	0.839	0.878	0.772	0.807	0.789	125	63	0.43	160.4
BPIC15 _{2f}	SM	0.783	0.877	0.828	0.514	0.596	0.552	129	49	0.36	0.6
	HPO _{sm}	0.808	0.851	0.829	0.561	0.582	0.572	133	56	0.30	1398.9
	RLS _{sm}	0.870	0.797	0.832	0.667	0.670	0.668	156	86	0.20	158.3
	ILS _{sm}	0.869	0.795	0.830	0.663	0.680	0.671	157	86	0.20	157.6
	TABU _{sm}	0.870	0.794	0.830	0.665	0.667	0.666	150	83	0.23	176.8
	SIMA _{sm}	0.871	0.775	0.820	0.677	0.662	0.669	159	93	0.26	167.4
BPIC15 _{3f}	SM	0.774	0.925	0.843	0.436	0.764	0.555	96	35	0.49	0.5
	HPO _{sm}	0.783	0.910	0.842	0.477	0.691	0.564	99	39	0.56	9230.4
	RLS _{sm}	0.812	0.903	0.855	0.504	0.775	0.611	110	53	0.35	151.5
	ILS _{sm}	0.833	0.868	0.850	0.533	0.775	0.631	120	66	0.23	153.8
	TABU _{sm}	0.832	0.852	0.842	0.558	0.690	0.617	121	64	0.23	173.4
	SIMA _{sm}	0.827	0.839	0.833	0.565	0.694	0.623	123	71	0.18	159.4
BPIC15 _{4f}	SM	0.762	0.886	0.820	0.516	0.615	0.562	101	37	0.27	0.5
	HPO _{sm}	0.785	0.860	0.821	0.558	0.578	0.568	103	40	0.27	736.4
	RLS _{sm}	0.825	0.854	0.839	0.634	0.672	0.652	114	57	0.21	146.9
	ILS _{sm}	0.853	0.807	0.829	0.649	0.657	0.653	117	64	0.27	147.8
	TABU _{sm}	0.811	0.794	0.803	0.642	0.661	0.651	115	61	0.24	161.7
	SIMA _{sm}	0.847	0.812	0.829	0.624	0.649	0.636	117	61	0.18	148.2
BPIC15 _{5f}	SM	0.806	0.915	0.857	0.555	0.598	0.576	110	38	0.34	0.6
	HPO _{sm}	0.789	0.941	0.858	0.529	0.655	0.585	102	30	0.33	972.3
	RLS _{sm}	0.868	0.813	0.840	0.737	0.731	0.734	137	78	0.14	159.3
	ILS _{sm}	0.868	0.813	0.840	0.737	0.731	0.734	137	78	0.14	153.8
	TABU _{sm}	0.885	0.818	0.850	0.739	0.746	0.743	137	79	0.14	173.3
	SIMA _{sm}	0.867	0.811	0.838	0.734	0.727	0.731	137	78	0.16	154.3

Table 2 continued

Event log	Discovery approach	Align. Acc.			Markov. Acc.			Complexity			Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.	
RTFMP	SM	0.996	0.958	0.977	0.959	0.311	0.470	22	17	0.46	2.9
	HPO _{sm}	0.887	1.000	0.940	0.685	0.696	0.690	20	9	0.35	2452.7
	RLS _{sm}	0.988	1.000	0.994	0.899	0.794	0.843	22	14	0.46	142.8
	ILS _{sm}	0.988	1.000	0.994	0.899	0.794	0.843	22	14	0.46	143.8
	TABU _{sm}	0.988	1.000	0.994	0.899	0.794	0.843	22	14	0.46	114.8
	SIMA _{sm}	0.986	1.000	0.993	0.875	0.893	0.884	23	15	0.39	131.0
SEPSIS	SM	0.764	0.706	0.734	0.349	0.484	0.406	32	23	0.94	0.4
	HPO _{sm}	0.925	0.588	0.719	0.755	0.293	0.423	33	34	0.39	28,846
	RLS _{sm}	0.839	0.630	0.720	0.508	0.430	0.466	35	29	0.77	145.4
	ILS _{sm}	0.812	0.625	0.706	0.455	0.436	0.445	35	28	0.86	157.1
	TABU _{sm}	0.839	0.630	0.720	0.508	0.430	0.466	35	29	0.77	137.0
	SIMA _{sm}	0.806	0.613	0.696	0.477	0.445	0.460	35	30	0.77	137.2

The values in bold are the best values of the measurements for each column within each block

**Fig. 8** BPIC14_f model discovered with SIMA_{sm} (above) and with SM (below)

time of SM (~ 1 second on average), it is much lower than HPO_{sm}, while consistently achieving higher accuracy.

4.3 Fodina

Tables 4 and 5 report the results of our comparative evaluation for Fodina. In these tables, we used “—” to report that a given accuracy measurement could not be reliably obtained due to

the unsoundness of the discovered process model. We held out from the tables two logs: BPIC12 and SEPSIS, because none of the six approaches (base APDA, hyper-parameter optimized and the four metaheuristics) was able to discover a sound process model. This is due to Fodina’s design which does not guarantee soundness.

Considering the remaining 18 logs, eleven times *all* the metaheuristics improved the Markovian F-score w.r.t. HPO_{f0}

Table 3 Comparative evaluation results for the proprietary logs—Split Miner

Event log	Discovery method	Align. Acc.			Markov. Acc.			(k = 5)			Complexity	Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.		
PRT1	SM	0.976	0.974	0.975	0.730	0.669	0.698	20	14	1.00	0.4	
	HPO _{sm}	0.999	0.948	0.972	0.989	0.620	0.762	19	14	0.53	298.3	
	RLS _{sm}	0.976	0.974	0.975	0.730	0.669	0.698	20	14	1.00	155.3	
	ILS _{sm}	0.976	0.974	0.975	0.730	0.669	0.698	20	14	1.00	153.2	
	TABU _{sm}	0.976	0.974	0.975	0.730	0.669	0.698	20	14	1.00	10.3	
	SIMA _{sm}	0.983	0.964	0.974	0.814	0.722	0.765	20	15	1.00	132.6	
PRT2	SM	0.795	0.581	0.671	0.457	0.913	0.609	29	23	1.00	0.3	
	HPO _{sm}	0.826	0.675	0.743	0.501	0.830	0.625	21	13	0.67	406.4	
	RLS _{sm}	0.886	0.421	0.571	0.629	0.751	0.685	29	34	1.00	141.4	
	ILS _{sm}	0.890	0.405	0.557	0.645	0.736	0.688	29	35	1.00	172.3	
	TABU _{sm}	0.866	0.425	0.570	0.600	0.782	0.679	29	33	1.00	143.1	
	SIMA _{sm}	0.886	0.424	0.574	0.629	0.751	0.685	29	34	1.00	139.7	
PRT3	SM	0.882	0.887	0.885	0.381	0.189	0.252	31	23	0.58	0.4	
	HPO _{sm}	0.890	0.899	0.895	0.461	0.518	0.488	26	14	0.81	290.2	
	RLS _{sm}	0.945	0.902	0.923	0.591	0.517	0.551	31	23	0.55	138.4	
	ILS _{sm}	0.945	0.902	0.923	0.591	0.517	0.551	31	23	0.55	144.2	
	TABU _{sm}	0.944	0.902	0.922	0.589	0.519	0.552	30	20	0.60	134.7	
	SIMA _{sm}	0.945	0.902	0.923	0.591	0.517	0.551	31	23	0.55	133.7	
PRT4	SM	0.884	1.000	0.938	0.483	1.000	0.652	25	15	0.96	0.5	
	HPO _{sm}	0.973	0.930	0.951	0.929	0.989	0.958	26	24	0.31	867.5	
	RLS _{sm}	0.997	0.903	0.948	0.993	0.990	0.992	26	28	0.92	140.1	
	ILS _{sm}	0.997	0.903	0.948	0.993	0.990	0.992	26	28	0.92	152.3	
	TABU _{sm}	0.955	0.914	0.934	0.883	0.988	0.932	26	26	0.77	138.6	
	SIMA _{sm}	0.997	0.903	0.948	0.993	0.990	0.992	26	28	0.92	136.9	
PRT6	SM	0.937	1.000	0.967	0.542	1.000	0.703	15	4	1.00	0.3	
	HPO _{sm}	0.937	1.000	0.967	0.542	1.000	0.703	15	4	1.00	105.1	
	RLS _{sm}	0.984	0.928	0.955	0.840	0.818	0.829	22	14	0.41	141.1	
	ILS _{sm}	0.984	0.928	0.955	0.840	0.818	0.829	22	14	0.41	144.2	
	TABU _{sm}	0.984	0.928	0.955	0.840	0.818	0.829	22	14	0.41	124.9	
	SIMA _{sm}	0.984	0.928	0.955	0.840	0.818	0.829	22	14	0.41	131.2	
PRT7	SM	0.914	0.999	0.954	0.650	1.000	0.788	29	10	0.48	0.6	
	HPO _{sm}	0.944	1.000	0.971	0.772	1.000	0.871	22	9	0.64	173.1	
	RLS _{sm}	0.993	1.000	0.996	0.933	1.000	0.965	23	11	0.78	139.2	
	ILS _{sm}	0.993	1.000	0.996	0.933	1.000	0.965	23	11	0.78	142.9	
	TABU _{sm}	0.993	1.000	0.996	0.933	1.000	0.965	23	11	0.78	134.0	
	SIMA _{sm}	0.993	1.000	0.996	0.933	1.000	0.965	23	11	0.78	131.9	
PRT10	SM	0.970	0.943	0.956	0.905	0.206	0.335	45	47	0.84	0.5	
	HPO _{sm}	0.936	0.943	0.939	0.810	0.243	0.374	30	22	0.73	1214.3	
	RLS _{sm}	0.917	0.989	0.952	0.741	0.305	0.432	44	41	0.86	153.0	
	ILS _{sm}	0.917	0.989	0.952	0.741	0.305	0.432	44	41	0.86	155.4	
	TABU _{sm}	0.917	0.989	0.952	0.741	0.305	0.432	44	41	0.86	117.6	
	SIMA _{sm}	0.917	0.989	0.952	0.741	0.305	0.432	44	41	0.86	136.7	

The values in bold are the best values of the measurements for each column within each block

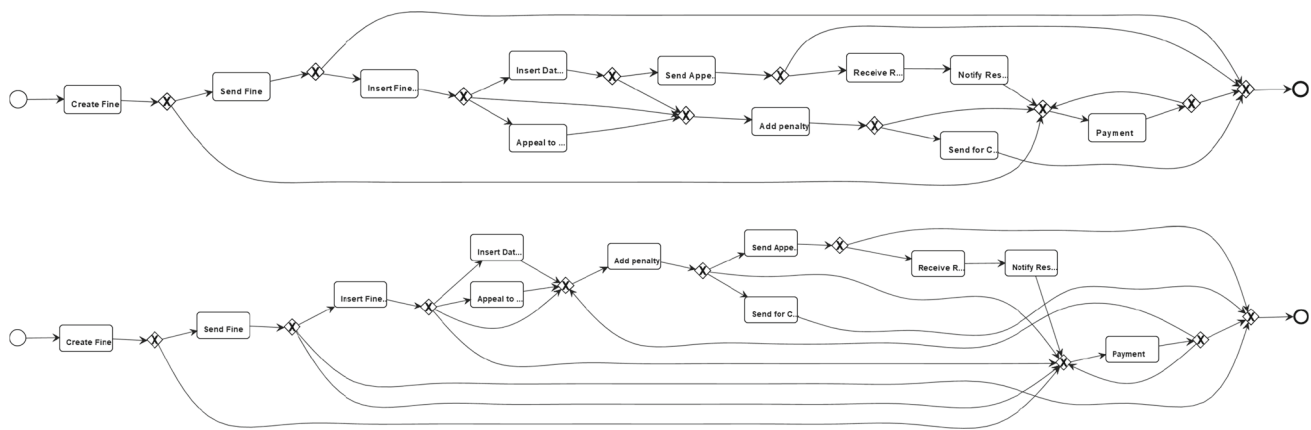


Fig. 9 RTFMP model discovered with SIMA_{sm} (above) and with SM (below)

Table 4 Comparative evaluation results for the public logs—Fodina

Event log	Discovery approach	Align. Acc.			Markov. Acc.			Complexity			Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.	
BPIC13 _{cp}	FO	0.999	0.879	0.935	0.997	0.647	0.784	13	10	0.77	0.1
	HPO _{fo}	0.999	0.879	0.935	0.997	0.647	0.784	13	10	0.77	17.7
	RLS _{fo}	0.994	0.963	0.978	0.947	0.864	0.904	12	9	0.67	290.6
	ILS _{fo}	0.994	0.880	0.934	0.935	0.758	0.837	12	8	0.92	151.2
	TABU _{fo}	0.994	0.963	0.978	0.947	0.864	0.904	12	9	0.67	95.2
	SIMA _{fo}	0.994	0.880	0.934	0.935	0.758	0.837	12	8	0.92	130.0
BPIC13 _{inc}	FO	0.994	0.877	0.932	0.950	0.576	0.717	13	10	0.85	0.291
	HPO _{fo}	0.994	0.877	0.932	0.950	0.576	0.717	13	10	0.85	112.0
	RLS _{fo}	0.994	0.877	0.932	0.950	0.576	0.717	13	10	0.85	304.7
	ILS _{fo}	0.994	0.877	0.932	0.950	0.576	0.717	13	10	0.85	180.1
	TABU _{fo}	0.998	0.743	0.852	0.987	0.604	0.749	14	15	1.00	129.0
	SIMA _{fo}	0.994	0.877	0.932	0.950	0.576	0.717	13	10	0.85	146.1
BPIC14 _f	FO	—	—	—	—	—	—	37	46	0.41	36.8
	HPO _{fo}	1.000	0.757	0.862	1.000	0.985	0.992	27	36	0.56	8612.7
	RLS _{fo}	1.000	0.757	0.862	1.000	0.985	0.992	27	36	0.56	370.7
	ILS _{fo}	1.000	0.757	0.862	1.000	0.985	0.992	27	36	0.56	365.5
	TABU _{fo}	1.000	0.757	0.862	1.000	0.985	0.992	27	36	0.56	358.5
	SIMA _{fo}	1.000	0.757	0.862	1.000	0.985	0.992	27	36	0.56	300.2
BPIC151 _f	FO	1.000	0.760	0.860	1.000	0.480	0.650	146	91	0.26	0.3
	HPO _{fo}	1.000	0.756	0.861	1.000	0.479	0.648	146	91	0.26	130.5
	RLS _{fo}	0.916	0.829	0.870	0.804	0.772	0.788	131	69	0.24	301.9
	ILS _{fo}	0.916	0.829	0.870	0.804	0.772	0.788	131	69	0.24	198.4
	TABU _{fo}	0.916	0.830	0.871	0.802	0.778	0.790	129	67	0.33	177.5
	SIMA _{fo}	0.918	0.833	0.873	0.777	0.799	0.788	127	67	0.34	174.4
BPIC152 _f	FO	—	—	—	—	—	—	195	159	0.09	48.5
	HPO _{fo}	—	—	—	—	—	—	187	145	0.11	118.7
	RLS _{fo}	—	—	—	—	—	—	181	131	0.09	306.0
	ILS _{fo}	—	—	—	—	—	—	175	120	0.11	276.1
	TABU _{fo}	0.876	0.754	0.810	0.653	0.608	0.630	177	120	0.13	262.3
	SIMA _{fo}	—	—	—	—	—	—	175	121	0.12	284.1

Table 4 continued

Event log	Discovery approach	Align. Acc.			Markov. Acc.			$(k = 5)$			Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.	
BPIC15 _{3f}	FO	—	—	—	—	—	—	174	164	0.06	4.3
	HPO _{fo}	0.983	0.601	0.746	0.925	0.208	0.339	163	161	0.07	402.9
	RLS _{fo}	—	—	—	—	—	—	166	141	0.07	303.5
	ILS _{fo}	0.924	0.713	0.805	0.701	0.444	0.543	158	131	0.10	247.1
	TABU _{fo}	—	—	—	—	—	—	163	131	0.09	235.5
	SIMA _{fo}	—	—	—	—	—	—	163	131	0.09	241.8
BPIC15 _{4f}	FO	—	—	—	—	—	—	157	127	0.15	1.3
	HPO _{fo}	0.995	0.660	0.793	0.973	0.302	0.461	153	126	0.14	443.0
	RLS _{fo}	0.887	0.790	0.836	0.708	0.610	0.655	127	77	0.17	308.3
	ILS _{fo}	0.882	0.801	0.839	0.697	0.628	0.661	127	75	0.17	300.5
	TABU _{fo}	0.864	0.806	0.834	0.675	0.652	0.663	127	74	0.17	274.5
	SIMA _{fo}	0.882	0.801	0.839	0.697	0.628	0.661	127	75	0.17	252.2
BPIC15 _{5f}	FO	1.000	0.698	0.822	1.000	0.362	0.532	166	125	0.15	2.4
	HPO _{fo}	1.000	0.698	0.822	1.000	0.362	0.532	166	125	0.15	238.1
	RLS _{fo}	0.886	0.810	0.846	0.727	0.703	0.715	150	94	0.11	303.4
	ILS _{fo}	0.884	0.819	0.850	0.719	0.724	0.722	147	90	0.13	268.1
	TABU _{fo}	0.886	0.814	0.849	0.723	0.730	0.727	149	92	0.11	217.1
	SIMA _{fo}	0.884	0.808	0.844	0.721	0.743	0.732	141	83	0.14	208.5
BPIC17 _f	FO	1.000	0.675	0.806	1.000	0.330	0.496	35	22	0.69	22.4
	HPO _{fo}	1.000	0.675	0.806	1.000	0.330	0.496	35	22	0.71	9755.7
	RLS _{fo}	0.999	0.675	0.806	0.997	0.331	0.497	33	20	0.70	309.9
	ILS _{fo}	0.999	0.675	0.806	0.997	0.331	0.497	33	20	0.70	313.5
	TABU _{fo}	0.999	0.675	0.806	0.997	0.331	0.497	33	20	0.70	305.7
	SIMA _{fo}	0.999	0.675	0.806	0.997	0.331	0.497	33	20	0.70	319.2
RTFMP	FO	0.996	0.933	0.964	0.937	0.148	0.256	31	32	0.19	0.4
	HPO _{fo}	0.884	1.000	0.939	0.646	0.857	0.737	18	7	0.56	2666.2
	RLS _{fo}	0.987	1.000	0.994	0.848	0.938	0.890	26	25	0.12	268.7
	ILS _{fo}	0.987	1.000	0.994	0.848	0.938	0.890	26	25	0.12	134.1
	TABU _{fo}	0.987	1.000	0.994	0.848	0.938	0.890	26	25	0.12	131.2
	SIMA _{fo}	0.987	1.000	0.993	0.847	0.923	0.883	28	27	0.11	133.9

The values in bold are the best values of the measurements for each column within each block

(and consequently FO), while 16 times at least one metaheuristic outperformed both FO and HPO_{fo}. The only two cases where none of the metaheuristics was able to discover a more accurate process model than HPO_{fo} were PRT2 and BPIC14_f. In the former log, this is because none of the metaheuristics discovered a sound process model within the given timeout of five minutes. However, we note that HPO_{fo} took almost four hours to discover a sound process model from the PRT2 log. In the latter log, this is because all the metaheuristics discovered the same model of HPO_{fo}.

Among the optimization metaheuristics, TABU_{fo} performed the best. This metaheuristic achieved 14 times out of 18 the highest Markovian F-score, followed by ILS (ten times). However, like for Split Miner, the differences in the

achieved F-score between the four metaheuristics are small. There is a difference of only 1–2 percentage points between the metaheuristics with highest F-score and the one with lowest F-score.

In the case of Fodina, the results achieved by the metaheuristics on the alignment-based F-score are more remarkable than the case of Split Miner, and in-line with the results obtained on the Markovian F-score. Indeed, 50% of the times, *all* the metaheuristics were able to outperform both FO and HPO_{fo} on the alignment-based F-score, and more than 80% of the times, at least one metaheuristic scored a higher alignment-based F-score than FO and HPO_{fo}. Such a result is remarkable considering that the objective function of the metaheuristics was the Markovian F-score.

Table 5 Comparative evaluation results for the proprietary logs—Fodina

Event log	Discovery method	Align. Acc.			Markov. Acc.			(k = 5)			Complexity	Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.		
PRT1	FO	—	—	—	—	—	—	30	28	0.53	0.2	
	HPO _{fo}	0.998	0.925	0.960	0.988	0.739	0.845	21	17	0.81	402.6	
	RLS _{fo}	0.988	0.964	0.976	0.888	0.827	0.857	21	16	0.86	302.7	
	ILS _{fo}	0.988	0.964	0.976	0.888	0.827	0.857	21	16	0.86	183.1	
	TABU _{fo}	0.994	0.957	0.976	0.981	0.844	0.907	21	17	0.86	149.3	
	SIMA _{fo}	0.988	0.964	0.976	0.888	0.827	0.857	21	16	0.86	154.0	
PRT2	FO	—	—	—	—	—	—	38	45	0.76	92.7	
	HPO _{fo}	1.000	0.276	0.432	0.998	0.148	0.258	29	78	1.00	12937.1	
	RLS _{fo}	—	—	—	—	—	—	48	56	0.08	301.0	
	ILS _{fo}	—	—	—	—	—	—	48	56	0.08	308.1	
	TABU _{fo}	—	—	—	—	—	—	53	70	0.08	313.0	
	SIMA _{fo}	—	—	—	—	—	—	—	—	—	854.9	
PRT3	FO	0.999	0.847	0.917	0.993	0.269	0.423	34	37	0.32	0.2	
	HPO _{fo}	—	—	—	—	—	—	73	93	0.18	756.5	
	RLS _{fo}	0.963	0.902	0.932	0.679	0.446	0.539	37	38	0.35	306.6	
	ILS _{fo}	0.963	0.902	0.932	0.679	0.446	0.539	37	38	0.35	157.3	
	TABU _{fo}	0.963	0.902	0.932	0.679	0.446	0.539	37	38	0.35	138.1	
	SIMA _{fo}	0.963	0.902	0.932	0.679	0.446	0.539	37	38	0.35	143.0	
PRT4	FO	—	—	—	—	—	—	37	40	0.54	46.0	
	HPO _{fo}	1.000	0.858	0.924	1.000	0.965	0.982	32	41	0.50	10914.5	
	RLS _{fo}	0.997	0.859	0.923	0.993	0.990	0.991	31	37	0.52	317.4	
	ILS _{fo}	0.997	0.903	0.948	0.993	0.993	0.993	27	32	0.74	314.4	
	TABU _{fo}	0.997	0.903	0.948	0.993	0.993	0.993	27	32	0.74	300.1	
	SIMA _{fo}	0.977	0.887	0.930	0.793	0.963	0.870	32	38	0.50	309.1	
PRT6	FO	1.000	0.908	0.952	1.000	0.632	0.775	22	17	0.41	0.1	
	HPO _{fo}	1.000	0.908	0.952	1.000	0.632	0.775	22	17	0.41	25.0	
	RLS _{fo}	0.984	0.928	0.955	0.840	0.818	0.829	22	14	0.41	278.8	
	ILS _{fo}	0.984	0.928	0.955	0.840	0.818	0.829	22	14	0.41	140.0	
	TABU _{fo}	0.984	0.928	0.955	0.840	0.818	0.829	22	14	0.41	129.3	
	SIMA _{fo}	0.984	0.928	0.955	0.840	0.818	0.829	22	14	0.41	131.2	
PRT7	FO	0.990	1.000	0.995	0.906	1.000	0.951	26	16	0.39	0.3	
	HPO _{fo}	0.990	1.000	0.995	0.906	1.000	0.951	26	16	0.39	50.2	
	RLS _{fo}	0.993	1.000	0.997	0.933	1.000	0.966	28	22	0.36	287.6	
	ILS _{fo}	0.993	1.000	0.997	0.933	1.000	0.966	28	22	0.36	140.3	
	TABU _{fo}	0.993	1.000	0.997	0.933	1.000	0.966	28	22	0.36	129.7	
	SIMA _{fo}	0.993	1.000	0.997	0.933	1.000	0.966	28	22	0.36	132.1	
PRT9	FO	—	—	—	—	—	—	32	45	0.72	53.1	
	HPO _{fo}	—	—	—	—	—	—	24	18	0.54	2799.5	
	RLS _{fo}	0.969	0.999	0.984	0.894	0.893	0.893	23	21	0.91	301.5	
	ILS _{fo}	—	—	—	—	—	—	34	26	0.15	263.2	
	TABU _{fo}	—	—	—	—	—	—	36	30	0.14	185.8	
	SIMA _{fo}	0.968	1.000	0.984	0.887	0.956	0.920	20	17	0.80	278.4	

Table 5 continued

Event log	Discovery method	Align. Acc.			Markov. Acc.			(k = 5)			Complexity	Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.		
PRT10	FO	0.990	0.922	0.955	0.961	0.087	0.159	52	85	0.64	0.2	
	HPO _{fo}	0.872	0.958	0.913	0.659	0.786	0.717	35	28	0.60	750.8	
	RLS _{fo}	0.964	0.965	0.965	0.870	0.813	0.840	44	46	0.25	301.1	
	ILS _{fo}	0.964	0.965	0.965	0.870	0.813	0.840	44	46	0.25	195.0	
	TABU _{fo}	0.964	0.965	0.965	0.870	0.813	0.840	44	46	0.25	165.0	
	SIMA _{fo}	0.965	0.963	0.964	0.874	0.809	0.840	44	47	0.25	161.2	

The values in bold are the best values of the measurements for each column within each block

Regarding the complexity of the models discovered by the metaheuristics, more than 50% of the times, it is lower than the complexity of the models discovered by FO and HPO_{fo}, and in the remaining cases in-line with the two baselines. Such a difference with the results we obtained for SM relates to the following two factors: (i) Split Miner discovers much simpler models than Fodina, and any further improvement is difficult to achieve; (ii) Fodina natively discovers more fitting models than Split Miner and hence, the metaheuristics aim at improving precision, ultimately removing model edges, and so reducing its complexity.

In terms of execution times, the four metaheuristics perform similarly, with an execution time between 150 and 300 s, slightly higher than the case of Split Miner.

4.4 Inductive Miner

Table 6 displays the results of our comparative evaluation for Inductive Miner. We held out from the table the five BPIC15 logs, because none of the three metaheuristics could discover a model within the five minutes timeout. This was due to scalability issues experienced by the Markovian accuracy, already known for the case of IM [6].

In the remaining 15 logs, 13 times *all* the metaheuristics improved the Markovian F-score w.r.t. IM_d, and only for the BPIC17_f log none of the metaheuristics could outperform IM_d. The best performing metaheuristic was SIMA_{imd}, eight times achieving the highest Markovian F-score, followed by TABU_{imd} and RLS_{imd}, which scored seven, respectively, six times the highest Markovian F-score. Again, we note that the differences in the achieved F-score across the four metaheuristics are small. There are several cases in which multiple metaheuristics achieve the same F-score, and a difference of only 1–2 percentage point between the best-performing and the worst-performing metaheuristics.

The results of the metaheuristics on the alignment-based F-score are similar to the case of Fodina, and they are broadly in-line with the results achieved on the Markovian F-score. Indeed, 80% of the times, *all* the metaheuristics were able to outperform IM_d, failing only in two logs out of 15.

Regarding the complexity of the models discovered by the metaheuristics, we recorded little variation w.r.t. the complexity of the models discovered by IM_d. Size and CFC did not notably improve nor worsen, except for the PRT9 and the BPIC14_f logs, where both size and CFC were reduced by about 30%.

In terms of execution times, the three metaheuristics perform similarly, with an average execution time close to 300 s, meaning that the majority of the times the solution-space exploration was interrupted by the timeout.

4.5 Discussion

The results of the evaluation show that the use of metaheuristics optimization brings consistent improvements in accuracy with respect to the baseline discovery approaches, in 80% of the cases. Furthermore, it produces consistently higher alignment-based F-score, even though this measure was not used as an objective function, due to the low scalability of alignment-based precision.

The drawback of using metaheuristics optimization is the longer execution time—several minutes versus less than a few seconds for the baselines.

In a small number of cases, the optimization framework did not yield any F-score improvement with respect to the corresponding unoptimized approach, due to: (i) a small solution-space (i.e., the baseline already discovers the best process model); or (ii) scalability issues (i.e., the Markovian accuracy could be computed within the timeout). While the former scenario is beyond our control and strictly relates to the complexity of the input event log, the latter reminds us of the limitations of the state-of-the-art accuracy measures (and especially precision) in the context of automated process discovery, and justifies our design choice of a modular optimization framework, that allows the use of new accuracy measures as objective functions in the future, which may be able to overcome such scalability issues.

Another remarkable finding is that the metaheuristically optimized versions of Split Miner and Fodina consistently outperform their hyper-parameter optimized counterparts.

Table 6 Comparative evaluation results for the public and proprietary logs—Inductive Miner

Event log	Discovery approach	Align. Acc.			Markov. Acc.			$(k = 5)$			Complexity	Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.		
BPIC12	IM _d	1.000	0.168	0.287	1.000	< 0.001	< 0.001	30	28	1.00	0.7	
	RLS _{imd}	0.661	0.763	0.708	0.220	0.163	0.187	40	21	1.00	300.0	
	TABU _{imd}	0.661	0.763	0.708	0.220	0.163	0.187	40	21	1.00	309.4	
	SIMA _{imd}	0.660	0.805	0.725	0.204	0.223	0.213	39	19	1.00	308.7	
BPIC13 _{cp}	IM _d	1.000	0.862	0.926	0.999	0.161	0.277	15	11	1.00	0.4	
	RLS _{imd}	0.984	0.889	0.934	0.882	0.424	0.573	9	5	1.00	301.4	
	TABU _{imd}	0.990	0.888	0.936	0.942	0.414	0.575	10	7	1.00	101.8	
	SIMA _{imd}	0.984	0.889	0.934	0.882	0.424	0.573	9	5	1.00	300.5	
BPIC13 _{inc}	IM _d	1.000	0.673	0.805	1.000	0.109	0.197	10	9	1.00	0.5	
	RLS _{imd}	0.895	0.921	0.908	0.679	0.517	0.587	10	6	1.00	301.5	
	TABU _{imd}	0.895	0.921	0.908	0.679	0.517	0.587	10	6	1.00	71.9	
	SIMA _{imd}	0.895	0.921	0.908	0.679	0.517	0.587	10	6	1.00	300.7	
BPIC14 _f	IM _d	0.861	0.782	0.820	0.507	0.814	0.625	27	16	1.00	0.8	
	RLS _{imd}	0.977	0.676	0.799	0.918	0.447	0.601	16	11	1.00	302.5	
	TABU _{imd}	0.949	0.673	0.788	0.859	0.505	0.636	17	13	1.00	303.3	
	SIMA _{imd}	0.977	0.676	0.799	0.918	0.447	0.601	16	11	1.00	300.8	
BPIC17 _f	IM _d	1.000	0.679	0.808	1.000	0.284	0.442	34	23	1.00	1.5	
	RLS _{imd}	0.674	0.815	0.738	0.241	0.214	0.227	27	11	1.00	302.7	
	TABU _{imd}	0.693	0.817	0.750	0.262	0.204	0.230	28	13	1.00	83.8	
	SIMA _{imd}	0.674	0.815	0.738	0.241	0.214	0.227	27	11	1.00	301.0	
RTFMP	IM _d	1.000	0.543	0.704	1.000	0.003	0.005	15	12	1.00	0.8	
	RLS _{imd}	0.938	0.886	0.911	0.784	0.379	0.511	21	14	1.00	321.1	
	TABU _{imd}	0.938	0.886	0.911	0.784	0.379	0.511	21	14	1.00	52.1	
	SIMA _{imd}	0.917	0.907	0.912	0.780	0.625	0.694	19	9	1.00	300.8	
SEPSIS	IM _d	1.000	0.291	0.451	0.918	0.006	0.012	24	23	1.00	0.4	
	RLS _{imd}	0.796	0.684	0.736	0.367	0.363	0.365	27	18	1.00	305.5	
	TABU _{imd}	0.796	0.684	0.736	0.367	0.363	0.365	27	18	1.00	306.9	
	SIMA _{imd}	0.813	0.581	0.678	0.482	0.310	0.377	25	16	1.00	301.6	
PRT1	IM _d	1.000	0.748	0.856	1.000	0.025	0.048	14	11	1.00	0.5	
	RLS _{imd}	0.974	0.946	0.960	0.692	0.707	0.699	16	10	1.00	304.4	
	TABU _{imd}	0.971	0.946	0.958	0.692	0.707	0.699	17	10	1.00	304.0	
	SIMA _{imd}	0.974	0.946	0.960	0.692	0.707	0.699	16	10	1.00	300.7	
PRT2	IM _d	1.000	0.243	0.390	1.000	0.109	0.196	13	11	1.00	0.9	
	RLS _{imd}	0.811	0.464	0.591	0.588	0.601	0.594	18	13	1.00	305.5	
	TABU _{imd}	0.788	0.461	0.581	0.542	0.566	0.554	16	11	1.00	303.0	
	SIMA _{imd}	0.792	0.413	0.543	0.524	0.674	0.590	18	13	1.00	307.3	
PRT3	IM _d	0.827	0.890	0.857	0.328	0.253	0.286	26	10	1.00	0.4	
	RLS _{imd}	0.914	0.896	0.905	0.501	0.593	0.543	26	14	1.00	305.1	
	TABU _{imd}	0.933	0.900	0.917	0.626	0.592	0.608	28	15	1.00	302.6	
	SIMA _{imd}	0.930	0.898	0.914	0.562	0.539	0.550	29	17	1.00	300.8	

Table 6 continued

Event log	Discovery approach	Align. Acc.			Markov. Acc.			$(k = 5)$			Complexity	Exec. time(s)
		Fitness	Prec.	F-score	Fitness	Prec.	F-score	Size	CFC	Struct.		
PRT4	IM _d	0.880	0.811	0.844	0.876	0.967	0.919	27	13	1.00	0.5	
	RLS _{imd}	0.962	0.879	0.919	1.000	0.956	0.977	19	13	1.00	301.0	
	TABU _{imd}	0.962	0.879	0.919	1.000	0.956	0.977	19	13	1.00	307.2	
	SIMA _{imd}	0.962	0.879	0.919	1.000	0.956	0.977	19	13	1.00	300.8	
PRT6	IM _d	0.917	0.988	0.951	0.524	0.350	0.420	18	6	1.00	0.4	
	RLS _{imd}	0.953	0.987	0.969	0.674	0.941	0.785	17	7	1.00	304.2	
	TABU _{imd}	0.905	0.915	0.910	0.488	0.903	0.634	18	10	1.00	643.0	
	SIMA _{imd}	0.953	0.987	0.969	0.674	0.941	0.785	17	7	1.00	300.6	
PRT7	IM _d	0.852	0.997	0.919	0.618	0.407	0.491	21	5	1.00	0.4	
	RLS _{imd}	0.917	1.000	0.957	0.700	1.000	0.824	20	6	1.00	305.5	
	TABU _{imd}	0.917	1.000	0.957	0.700	1.000	0.824	20	6	1.00	1013.9	
	SIMA _{imd}	0.960	0.988	0.974	0.664	0.752	0.705	23	13	1.00	309.6	
PRT9	IM _d	0.586	0.461	0.516	0.078	0.014	0.024	22	15	1.00	2.6	
	RLS _{imd}	0.945	1.000	0.972	0.851	1.000	0.919	16	9	1.00	304.2	
	TABU _{imd}	0.946	1.000	0.972	0.856	0.947	0.899	18	10	1.00	306.6	
	SIMA _{imd}	0.954	1.000	0.976	0.890	0.909	0.899	15	8	1.00	300.0	
PRT10	IM _d	0.530	0.656	0.586	0.386	0.000	0.001	36	28	1.00	0.5	
	RLS _{imd}	0.859	0.961	0.907	0.664	0.691	0.677	30	24	1.00	300.3	
	TABU _{imd}	0.912	0.907	0.909	0.790	0.484	0.600	30	24	1.00	33.1	
	SIMA _{imd}	0.862	0.941	0.900	0.671	0.719	0.694	32	28	1.00	307.6	

The values in bold are the best values of the measurements for each column within each block

This means that the space of possible process models that can be explored by tweaking the parameters in input (e.g., the noise filter threshold) is not as rich as the space of process models that can be generated by repeatedly perturbing the DFG.

Finally, we found that all four metaheuristics considered in the evaluation led to similar F-scores. The differences in F-score between the best-performing and the worst-performing metaheuristics are generally negligible, in the order of 1–2 percentage points. For Inductive Miner, all four metaheuristics end up exploring the search space in a similar manner, leading to the same results in several cases. This may be explained by the fact that the set of possible models that Inductive Miner can generate is narrower than that generated by Fodina or Split Miner, because Inductive Miner can only generate block-structured models.

5 Conclusion

This paper showed that the use of S-metaheuristics is a promising approach to enhance the accuracy of DFG-based automated process discovery approaches. The outlined approach takes advantage of the DFG's simplicity to define

efficient perturbation functions that improve fitness or precision while preserving structural properties required to ensure model correctness.

The evaluation showed that the metaheuristically optimized approaches consistently achieve higher accuracy than the corresponding unoptimized baselines (Split Miner, Fodina, and Inductive Miner—directly-follows). This observation holds both when measuring accuracy via Markovian F-score and via alignment-based F-score. As expected, these accuracy gains come at the expense of higher execution times. This is naturally given that the metaheuristics needs to execute the baseline approach several hundred times and it needs to measure the accuracy of each of the resulting process models. The evaluation also showed that the choice of optimization metaheuristic (among those considered in the paper) does not have a substantial effect on the accuracy of the resulting process models, nor on their size or complexity.

In its current form, the framework focuses on improving F-score. In principle, the framework could also be used to optimize other objective functions, such as model complexity, measured by number of edges or control-flow complexity measures, for example. Related to the above, the framework could be extended to optimize multiple dimensions simultaneously, using multi-objective (Pareto-front) optimization

techniques instead of single-objective ones. Along the same lines, it may also be possible to adapt the framework in order to optimize one measure (e.g., F-score), subject to one or more constraints on other measures (e.g., that the number of edges in the discovered model must be below a given threshold). Lastly, another opportunity for future work may be the automation of the best metaheuristic selection, without compromising the time performance of the framework.

Another limitation of the proposed approach is that the DFG perturbations employed in the optimization phase do not use the frequencies of the directly-follows relations (i.e., arc frequencies in the DFG are not used). In other words, the proposed approach makes the following two design choices: (i) the perturbations either add an arc or remove an arc but they do not alter the frequency of an arc; and (ii) the decision as to which arcs to add or remove is not taken based on arc frequencies. The rationale for the first of these design choices is that modifying the arc frequencies would require us to have a criterion for deciding by how much should be frequencies be altered. This criterion would have to be dependent on the way the underlying automated process discovery algorithm uses the arc frequencies. We opted not to do in order to ensure the optimization method is independent of the underlying base algorithm. The rationale for the second choice is that the perturbations should have an element of randomness in order to allow the metaheuristics to explore a wider subset of the search space. Poor perturbations are likely to lead to solutions with lower F-scores, which are eventually discarded by the metaheuristics, but a transformation that leads to solutions with lower F-scores may later give rise to other solutions with higher F-score as the search unfolds. This having been said, it is possible that perturbations that remove arcs based on arc frequency might help the heuristics to focus on areas of the search space with higher F-score. A direction for future work is to explore other perturbation heuristics including frequency-aware ones.

A third limitation of the framework is that it only considers four S-metaheuristics. There is room for investigating further metaheuristics such as variants of simulated annealing, e.g., using different cooling schedules. Along a similar direction, this study could be extended to investigate the trade-offs between S-metaheuristics and P-metaheuristics in this setting.

Finally, the evaluation put into evidence scalability limitations of the Markovian precision measure for some datasets. These limitations are not specific to this precision measure—they also apply, and sometimes to a larger extent, to other precision measures including ETC precision and entropy-based precision [26]. There is a need for more scalable measures of precision in order to make metaheuristic optimization more broadly applicable in the context of automated process discovery.

Acknowledgements We thank Raffaele Conforti for his input to an earlier version of this paper. This research is partly funded by the Australian Research Council (DP180102839) and the European Research Council (PIX Project).

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Adriansyah, A., Munoz-Gama, J., Carmona, J., van Dongen, B., van der Aalst, W.: Measuring precision of modeled behavior. *ISEB* **13**(1), 37–67 (2015)
- Adriansyah, A., van Dongen, B., van der Aalst, W.: Conformance checking using cost-based fitness analysis. In: *EDOC. IEEE* (2011)
- Alizadeh, S., Norani, A.: ICMA: a new efficient algorithm for process model discovery. *Appl. Intell.* **48**(11), 4497–4514 (2018)
- Atashpaz-Gargari, E., Lucas, C.: Imperialist competitive algorithm: an algorithm for optimization inspired by imperialistic competition. In: *Evolutionary computation, 2007. CEC 2007. IEEE Congress on*, pp 4661–4667. *IEEE* (2007)
- Augusto, A., Armas Cervantes, A., Conforti, R., Dumas, M., La Rosa, M.: Measuring fitness and precision of automatically discovered process models: a principled and scalable approach. *IEEE Trans. Knowl. Data Eng.* (2020). To appear <https://doi.org/10.1109/TKDE.2020.3003258>
- Augusto, A., Armas-Cervantes, A., Conforti, R., Dumas, M., La Rosa, M., Reissner, D.: Abstract-and-compare: a family of scalable precision measures for automated process discovery. In: *BPM. Springer, Berlin* (2018)
- Augusto, A., Conforti, R., Dumas, M., La Rosa, M., Bruno, G.: Automated discovery of structured process models from event logs: the discover-and-structure approach. *DKE* **117**, 373–392 (2017)
- Augusto, A., Conforti, R., Dumas, M., La Rosa, M., Maggi, F.M., Marrella, A., Mecella, M., Soo, A.: Automated discovery of process models from event logs: Review and benchmark. *IEEE Trans. Knowl. Data Eng.* **31**(4), 686–705 (2019)
- Augusto, A., Conforti, R., Dumas, M., La Rosa, M., Polyvyanny, A.: Split miner: automated discovery of accurate and simple business process models from event logs. *KAIS* **59**(2), 251–284 (2018)
- Augusto, A., Dumas, M., La Rosa, M.: Metaheuristic optimization for automated business process discovery. In: *BPM. Springer, Berlin* (2019)
- Bäck, Thomas, Schwefel, Hans-Paul: An overview of evolutionary algorithms for parameter optimization. *Evol. Comput.* **1**(1), 1–23 (1993)
- Boussaïd, I., Lepagnot, J., Siarry, P.: A survey on optimization metaheuristics. *Inf. Sci.* **237**, 82–117 (2013)
- Buijs, J., van Dongen, B., van der Aalst, W.: On the role of fitness, precision, generalization and simplicity in process discovery. In: *CoopIS. Springer, Berlin* (2012)
- Burattin, A., Sperduti, A.: Automatic determination of parameters' values for heuristics miner++. In: *IEEE Congress on Evolutionary Computation* (2010)

15. Chifu, V.R., Pop, C.B., Salomie, I., Balla, I., Paven, R.: Hybrid particle swarm optimization method for process mining. In: ICCP. IEEE (2012)
16. de Medeiros, A.K.: Genetic process mining. PhD thesis, Eindhoven University of Technology (2006)
17. Effendi, Y.A., Sarno, R.: Discovering optimized process model using rule discovery hybrid particle swarm optimization. In: 2017 3rd International Conference on Science in Information Technology (ICSITech), pp. 97–103. IEEE (2017)
18. Gao, D., Liu, Q.: An improved simulated annealing algorithm for process mining. In: CSCWD. IEEE (2009)
19. Glover, F.: Future paths for integer programming and links to artificial intelligence. *Comput Oper. Res.* **13**(5), 533–549 (1986)
20. Holland, J.H.: *Adaptation in Natural and Artificial Systems: An Introductory Analysis with Applications to Biology, Control, and Artificial Intelligence*. MIT Press, Cambridge (1992)
21. Kennedy, J.: Particle swarm optimization. In: *Encyclopedia of Machine Learning*, pp. 760–766. Springer, Berlin (2011)
22. Kirkpatrick, S., Gelatt, C.D., Vecchi, M.P.: Optimization by simulated annealing. *Science* **220**, 4598 (1983)
23. Kuhn, H.W.: The Hungarian method for the assignment problem. *NRL* **2**(1–2), 83–97 (1955)
24. Leemans, S.J.J., Fahland, D., van der Aalst, W.M.P.: Scalable process discovery and conformance checking. *Softw. Syst. Model.* **17**(2), 599–631 (2018)
25. Nurlaili, A.L., Sarno, R.: A combination of the evolutionary tree miner and simulated annealing. In: 2017 4th International Conference on Electrical Engineering, Computer Science and Informatics (EECSI), pp. 1–5. IEEE (2017)
26. Polyvyanyy, A., Solti, A., Weidlich, M., Ciccio, C.D., Mendling, J.: Monotone precision and recall measures for comparing executions and specifications of dynamic systems. *CoRR* (2018). [arXiv:1812.07334](https://arxiv.org/abs/1812.07334)
27. Ribeiro, J., Carmona Vargas, J.: A method for assessing parameter impact on control-flow discovery algorithms. In: *Algorithms & Theories for the Analysis of Event Data* (2015)
28. Smirnov, S., Weidlich, M., Mendling, J.: Business process model abstraction based on behavioral profiles. In: *International Conference on Service-Oriented Computing*, pp 1–16. Springer, Berlin (2010)
29. Song, W., Liu, S., Liu, Q.: Business process mining based on simulated annealing. In: *ICYCS*. IEEE (2008)
30. Stützle, T.: Local search algorithms for combinatorial problems. PhD thesis, Darmstadt University of Technology (1998)
31. Suriadi, S., Andrews, R., ter Hofstede, A.H.M., Wynn, M.T.: Event log imperfection patterns for process mining: towards a systematic approach to cleaning event logs. *Inf. Syst.* **64**, 132–150 (2017)
32. Syring, A.F., Tax, N., van der Aalst, W.M.P.: Evaluating conformance measures in process mining using conformance propositions. In: *Transactions on Petri Nets and Other Models of Concurrency XIV*, pp. 192–221. Springer, Berlin (2019)
33. van der Aalst, W.: *Process Mining: Data Science in Action*. Springer, Berlin (2016)
34. vanden Broucke, S., De Weerd, J.: Fodina: a robust and flexible heuristic process discovery technique. *DSS* **100**, 109–118 (2017)
35. Weijters, A., Ribeiro, J.: Flexible heuristics miner (FHM). In: *CIDM*. IEEE (2011)
36. Wen, L., Van Der Aalst, W.M.P., Wang, J., Sun, J.: Mining process models with non-free-choice constructs. *Data Min. Knowl. Discov.* **15**(2), 145–180 (2007)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Adriano Augusto is an Associate Lecturer at the University of Melbourne. He completed his joint Ph.D. at the University of Tartu and the University of Melbourne, in 2020, by defending a thesis focused on automated process discovery. His teaching and research interests include business process management, process mining, automated process discovery, and robotic process automation.



Marlon Dumas is Professor of Information Systems at University of Tartu, Estonia. His research interests span across the fields of software engineering, information systems and business process management. He is co-author of the textbook "Fundamentals of Business Process Management" (Springer, 2nd edition).

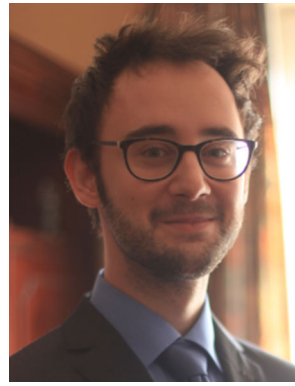


Marcello La Rosa is Professor of Information Systems at the University of Melbourne, Australia. His research interests include process mining, consolidation and automation. He leads the Aprimore Initiative, a cross-university collaboration for the development of an advanced process analytics platform, and co-authored the textbook "Fundamentals of Business Process Management" (Springer, 2nd edition).



Sander J.J. Leemans is an Assistant Professor (Lecturer) at the Queensland University of Technology, Australia, in the school of Information Systems. His research interests include process mining, process discovery, conformance checking, stochastic process mining, and robotic process automation. In particular, he specialises in making solid academic techniques available to end-users, analysts, and industry partners. He teaches business process management, business process modeling

and business process improvement.



Seppe vanden Broucke received a PhD in Applied Economics at KU Leuven, Belgium in 2014. Currently, Seppe is working as an assistant professor at the department of Business Informatics at UGent (Belgium) and is a lecturer at KU Leuven (Belgium). Seppe's research interests include business data mining and analytics, machine learning, process management, process mining. His work has been published in well-known international journals and presented at top conferences.