

## Genetic process mining: an experimental evaluation

A. K. A. de Medeiros · A. J. M. M. Weijters ·  
W. M. P. van der Aalst

Received: 13 March 2006 / Accepted: 20 November 2006 / Published online: 31 January 2007  
Springer Science+Business Media, LLC 2007

**Abstract** One of the aims of process mining is to retrieve a process model from an event log. The discovered models can be used as *objective* starting points during the deployment of process-aware information systems (Dumas et al., eds., *Process-Aware Information Systems: Bridging People and Software Through Process Technology*. Wiley, New York, 2005) and/or as a feedback mechanism to check prescribed models against enacted ones. However, current techniques have problems when mining processes that contain non-trivial constructs and/or when dealing with the presence of noise in the logs. Most of the problems happen because many current techniques are based on *local* information in the event log. To overcome these problems, we try to use genetic algorithms to mine process models. The main motivation is to benefit from the *global search* performed by this kind of algorithms. The non-trivial constructs are tackled by choosing an internal representation that supports them. The problem of noise is naturally tackled by the genetic algorithm because, per definition, these algorithms are robust to noise. The main challenge in a genetic approach is the definition of a good fitness measure because it guides the global search performed by the genetic algorithm. This paper explains how the genetic algorithm works. Experiments with synthetic and real-life logs show that the fitness measure indeed leads to the mining of process models that are *complete* (can reproduce all the behavior in the log) and *precise* (do not allow for extra behavior that cannot be derived from the event log). The genetic algorithm is implemented as a plug-in in the ProM framework.

---

Responsible editor: Eamonn Keogh.

---

A. K. A. de Medeiros (✉) · A. J. M. M. Weijters · W. M. P. van der Aalst  
Department of Technology Management, Eindhoven University of Technology, P.O. Box 513,  
5600 MB Eindhoven, The Netherlands.  
e-mail: a.k.medeiros@tm.tue.nl

**Keywords** Process mining · Genetic mining · Genetic algorithms · Petri nets · Workflow nets

## 1 Introduction

Today's organizations are supported by a wide variety of information systems. Some systems only support a single task (e.g., a text editor). However, most organizations are using systems that support processes, i.e., not a single task but the glue between tasks. Examples are WorkFlow Management (WFM) systems and Enterprise Resource Planning (ERP) systems. Typically, these systems record events that can be linked to the execution of some task in the process. Therefore, it makes sense to analyze these events to get feedback about enacted processes. Buzzwords such as Business Process Intelligence (BPI) and Business Activity Monitoring (BAM) indicate the interest of organizations and software developers in solutions able to extract knowledge from so-called event logs. However, most of the commercial systems (e.g., Cognos and Business Objects) focus on exclusively performance issues such as flow time and utilization. These systems abstract from the process itself and can only be applied if the process is well-defined and fixed. ARIS PPM is one of the few commercial systems actually trying to discover more information by monitoring events. One of the reasons for this limited support is that it is very difficult to extract process knowledge without having some a-priori process model. This triggered the development of process mining techniques that aim at automatically discovering process models based on event logs.

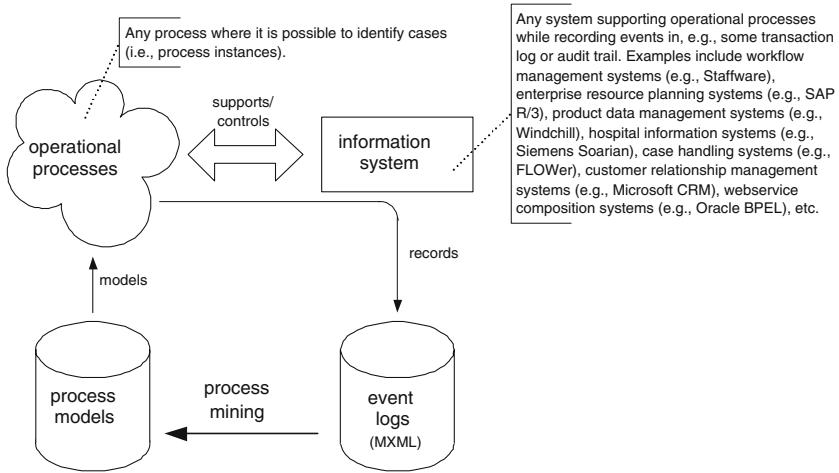
Figure 1 shows the concept of process mining. Some operational process is supported by some information system that records events in some event log. This event log is used to extract process models that describe the observed behavior. This information is valuable to better understand processes and to improve them. In our experience, real processes tend to deviate from the idealistic processes people have in mind. The practical relevance of process mining is obvious. Unfortunately, existing techniques have severe limitations.<sup>1</sup> Therefore, we present a new approach using a genetic algorithm. However, before we introduce our approach, we first need to clarify the concept of process mining.

### 1.1 Process mining

One of the aims of process mining is to *automatically* build a process model that describes the behavior contained in an event log. The models mined by process mining tools can be used as an *objective starting point* during the deployment of systems that support the execution of processes and/or as a *feedback mechanism* to check the prescribed process model against the enacted one. We use

---

<sup>1</sup> In Sect. 1.2, we elaborate more on these limitations.

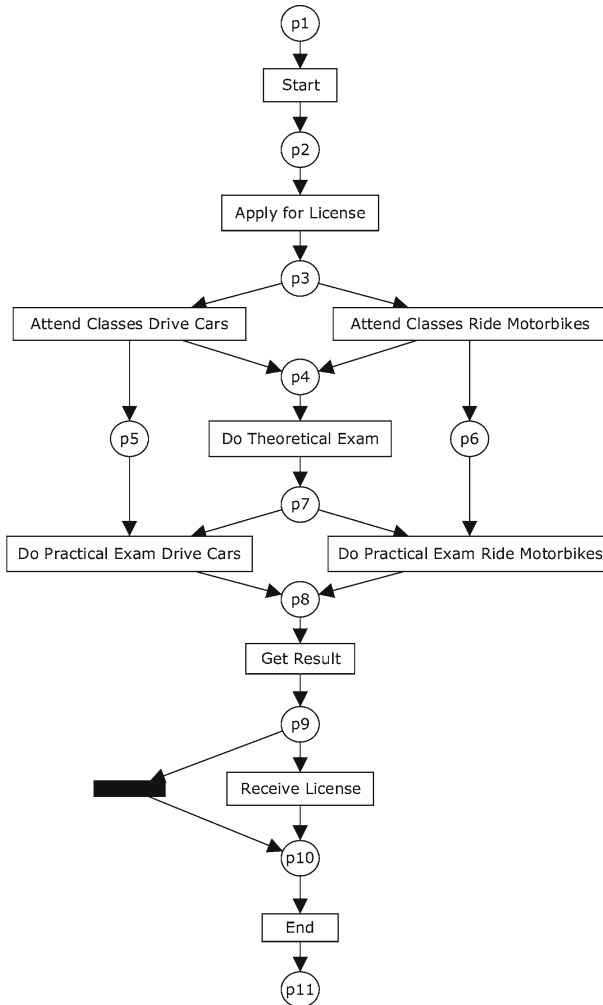


**Fig. 1** Overview of process mining

**Table 1** Example of an event log with four process instances

Identifier	Process instance
1	Start, Apply for License, Attend Classes Drive Cars, Do Theoretical Exam, Do Practical Exam Drive Cars, Get Result, End
2	Start, Apply for License, Attend Classes Ride Motorbikes, Do Theoretical Exam, Do Practical Exam Ride Motorbikes, Get Result, Receive License, End
3	Start, Apply for License, Attend Classes Drive Cars, Do Theoretical Exam, Do Practical Exam Drive Cars, Get Result, Receive License, End
4	Start, Apply for License, Attend Classes Ride Motorbikes, Do Theoretical Exam, Do Practical Exam Ride Motorbikes, Get Result, End

an example to illustrate how process mining techniques work. Consider the event log shown in Table 1. This log shows the event traces (process instances) for four different applications to get a license to ride motorbikes or drive cars. Note that applicants for different types of licenses do the same theoretical exam (task “Do Theoretical Exam”) but different practical ones (tasks “Do Practical Exam Drive Cars” or “Do Practical Exam Ride Motobikes”). In other words, whenever the task “Attend classes Drive Cars” is executed, the task “Do practical Exam Drive Cars” is the only one that can be executed after the applicant has done the theoretical exam. This shows that there is a *non-local* dependency between the tasks “Attend Classes Drive Cars” and “Do Practical Exam Drive Cars,” and also between the tasks “Attend Classes Ride Motorbikes” and “Do Practical Exam Ride Motorbikes.” The dependency is non-local because it cannot be detected by simply looking at the direct predecessor and



**Fig. 2** Mined net for the log in Table 1

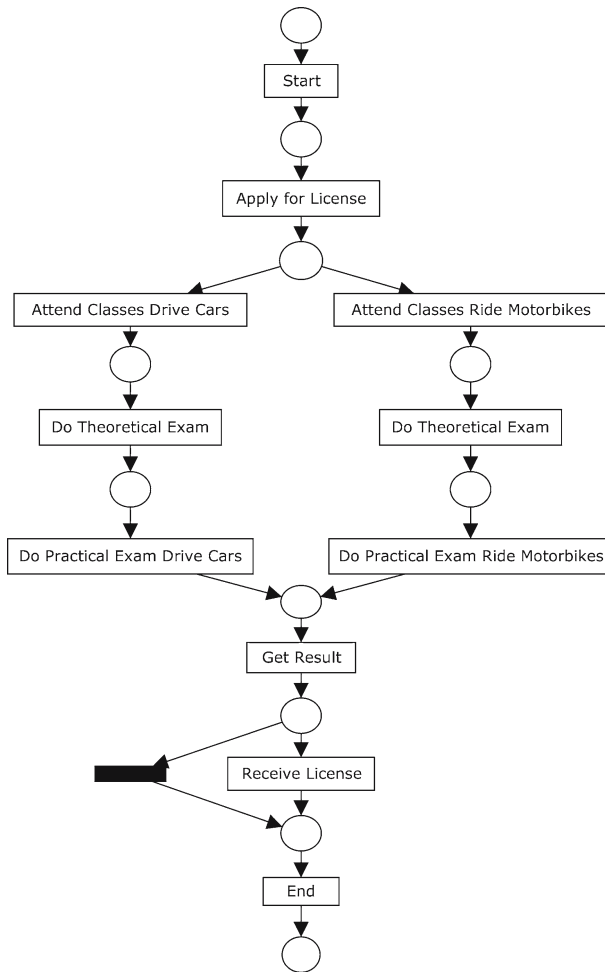
successor of those tasks in the log in Table 1, e.g., “Attend Classes Drive Cars” is never followed directly by “Do Practical Exam Drive Cars”. Moreover, note that only in some process instances (two and three) the task “Receive License” was executed. These process instances point to the cases in which the candidate passed the exams. Based on this log and these observations, process mining tools could be used to retrieve the model in Fig. 2. In this case, we are using Petri nets (Desel and Esparza 1995; Murata 1989) to depict this model. We do so because Petri nets will be used to explain the semantics of our internal representation. Moreover, we use Petri net-based analysis techniques to analyze the resulting models. Using the Petri net representation, our tools allow for the automatic translation of the discovered model to a variety of modelling

notations including Event-Driven Process Chains (used by ARIS, ARIS PPM, SAP) and YAWL (an open source workflow system).

Petri nets are a formalism to model concurrent processes. Graphically, Petri nets are bipartite directed graphs with two node types: *places* and *transitions*. Places represent conditions in the process. Transitions represent actions. Tasks in the event logs correspond to transitions in Petri nets. The state of a Petri net (or process for us) is described by adding tokens (black dots) to places. The dynamics of the Petri net is determined by the *firing rule*. A transition can be executed (i.e., an action can take place in the process) when all of its input places (i.e., pre-conditions) have at least a number of tokens that is equal to the number of directed arcs from the place to the transition. After execution, the transition removes tokens from the input places (one token is removed for every input arc from the place to the transition) and produces tokens for the output places (again, one token is produced for every output arc). Besides, the Petri nets that we consider have a single *start* place and a single *end* place. This means that the processes we describe have a single start point and a single end point. For the Petri net in Fig. 2, in the initial state there is only one token in place “p1.” This implies that “Start” is the only transition that can be executed in the initial state. When “Start” executes (or fires), one token is removed from the place “p1” and one token is added to the place “p2.” In a similar way, the firing of “Apply for License” marks place “p3.” In this marking, “Attend Classes Drive Cars” or “Attend Classes Ride Motorbikes” can fire. If “Attend Classes Drive Cars” fires, it consumes the token in “p3” and produces one token for “p4” and another for “p5.” Note that, although the place “p5” has now one token, the transition “Do Practical Exam Drive Cars” cannot fire yet because the place “p7” is not marked. The enabling and firing of transitions proceeds in a similar way until the place “p11” is marked.

## 1.2 Limitations of current approaches

Current research in process mining (van der Aalst et al. 2003, 2004; Agrawal et al. 1998; Greco et al. 2005; Herbst and Karagiannis 2004; Cook et al. 2004; Alves de Medeiros et al. 2003, 2004a; Schimm 2004; Wen et al. 2006) still has problems to discover process models with certain structural constructs and/or to deal with the presence of noise in the logs (cf. Sect. 9). The main problematic constructs are: non-free-choice, invisible tasks and duplicate tasks (Alves de Medeiros et al. 2003). *Non-free-choice* constructs combine synchronization and choice. The example in Fig. 2 shows a non-free-choice construct involving the tasks “Do Practical Exam Drive Cars” and “Do Practical Exam Ride Motorbikes.” The current techniques do not capture the dependency between (1) the tasks “Attend Classes Drive Cars” and “Do Practical Exam Drive Cars,” and (2) the tasks “Attend Classes Ride Motorbikes” and “Do Practical Exam Ride Motorbikes.” *Invisible tasks* are only used for routing purposes and do not appear in the log. For instance, the process in Fig. 2 has an invisible task to skip the execution of the task “Receive License.” Current techniques have



**Fig. 3** Another model that correctly portrays the behavior in the log in Table 1. Note that this model uses *duplicate tasks* instead of the *non-free-choice construct* in Fig. 2

difficulties discovering these routing tasks because they do not appear in the log. *Duplicate tasks* means that multiple transitions have the same label in the original process model. The problem here is that most of the mining techniques treat these duplicate tasks as a single one. For instance, Fig. 3 shows a model that also captures the behavior in the log in Table 1 by duplicating the task “Do Theoretical Exam.” *Noise* characterizes low-frequent behavior in the log. It can appear in two situations: event traces were somehow incorrectly logged (for instance, due to temporary system misconfiguration) or event traces reflect exceptional situations. Either way, most of the techniques will try to find a process model that can parse all the traces in the log. However, the presence of noise may hinder the correct mining of the most common behavior.

One of the reasons why the current techniques typically cannot cope with the above mentioned problematic constructs and/or with noisy logs is because their search is based on *local* information in the log. For instance, the  $\alpha$ -algorithm (see van der Aalst et al. 2004 for details) uses only information about which tasks directly succeed or precede one another in the process instances. As a result, this algorithm does not capture the dependency in non-free-choice constructs. For example, the  $\alpha$ -algorithm will never discover the Petri net in Fig. 2, for the log in Table 1, because none of the process instances has the sub-trace “Attend Classes Drive Cars, Do Practical Exam Drive Cars” or “Attend Classes Ride Motorbikes, Do Practical Exam Ride Motorbikes.” Consequently, the  $\alpha$ -algorithm will not link these tasks.

### 1.3 Genetic process mining

To overcome the limitations of the current process mining techniques, our research uses *genetic algorithms* (Eiben and Smith 2003) to mine process models. The main motivation is to benefit from the *global* search that is performed by this kind of algorithms.

Genetic algorithms are adaptive search methods that try to mimic the process of evolution. These algorithms start with an initial population of individuals. Every individual is assigned a fitness measure to indicate its quality. In our case, an individual is a possible process model and the fitness is a function that evaluates how well the individual is able to reproduce the behavior in the log. Populations evolve by selecting the fittest individuals and generating new individuals using genetic operators such as *crossover* (combining parts of two or more individuals) and *mutation* (random modification of an individual).

When using genetic algorithms to mine process models, there are three main concerns. The first is to define the *internal representation*. The internal representation defines the search space of a genetic algorithm. The internal representation that we define and explain in this paper supports all the problematic constructs, except for duplicate tasks. The second concern is to define the *fitness measure*. In our case, the fitness measure evaluates the quality of a point (individual or process model) in the search space against the event log. A genetic algorithm searches for individuals whose fitness is maximal. Thus, our fitness measure makes sure that individuals with a maximal fitness can parse all the process instances (traces) in the log and, ideally, not more than those traces. The reason for this is that we aim at discovering a process model that reflects as close as possible the behavior expressed in the event log. If the mined model allows for lots of extra behavior that cannot be derived from the log, it does not give a precise description of what is actually happening. The third concern relates to the *genetic operators* (crossover and mutation) because they should ensure that all points in the search space defined by the internal representation may be reached when the genetic algorithm runs. This paper presents a genetic algorithm that addresses these three concerns.

## 1.4 Road map

The remainder of the paper is organized as follows. Section 2 introduces the main definitions of Petri nets that are used in this paper. Section 3 explains the internal representation that we use and defines its semantics by mapping it onto Petri nets. Section 4 presents the genetic algorithm to mine processes that may have arbitrary mixtures of choice and synchronization (i.e., non-free-choice constructs) and invisible tasks. The genetic algorithm is also robust to noisy logs. Section 5 explains the metrics that we have developed to assess the quality of mined models while conducting the experiments. Section 6 discusses the experiments and results. The experiments include synthetic logs. Section 7 shows the results of applying the genetic algorithm to logs from a municipality in The Netherlands. Section 8 compares the results of the genetic algorithm with the results obtained by two other related process mining techniques. Section 9 discusses the related work. Section 10 contains the conclusions and future work.

## 2 Preliminaries

This section introduces standard Petri-net notations that are used to explain the semantics of the internal representation of our genetic algorithm.

### 2.1 Petri nets

We use a variant of the classic Petri-net model, namely *Place/Transition nets*. For an elaborate introduction to Petri nets, the reader is referred to [Desel and Esparza \(1995\)](#), [Murata \(1989\)](#), and [Reisig and Rozenberg \(1998\)](#).

**Definition 1** (*P/T-nets*)<sup>2</sup> A Place/Transition net, or simply P/T-net, is a tuple  $(P, T, F)$  where:

- (1)  $P$  is a finite set of places,
- (2)  $T$  is a finite set of transitions such that  $P \cap T = \emptyset$ , and
- (3)  $F \subseteq (P \times T) \cup (T \times P)$  is a set of directed arcs, called the flow relation.

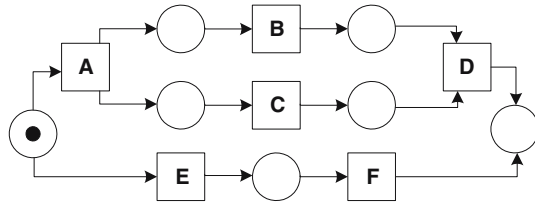
A marked P/T-net is a pair  $(N, s)$ , where  $N = (P, T, F)$  is a P/T-net and where  $s$  is a bag over  $P$  denoting the *marking* of the net, i.e.,  $s \in P \rightarrow \mathbb{N}$ . The set of all marked P/T-nets is denoted  $\mathcal{N}$ .

A marking is a *bag* over the set of places  $P$ , i.e., it is a function from  $P$  to the natural numbers. We use square brackets for the enumeration of a bag, e.g.,  $[a^2, b, c^3]$  denotes the bag with two  $a$ s, one  $b$ , and three  $c$ s. The sum of two bags  $(X + Y)$ , the difference  $(X - Y)$ , the presence of an element in a bag  $(a \in X)$ ,

<sup>2</sup> In the literature, the class of Petri nets introduced in Definition 1 is sometimes referred to as the class of (unlabeled) *ordinary* P/T-nets to distinguish it from the class of Petri nets that allows more than one arc between a place and a transition, and the class of Petri nets that allows for transition labels.



**Fig. 4** An example of a Place/Transition net



the intersection of two bags  $(X \cap Y)$  and the notion of subbags  $(X \leq Y)$  are defined in a straightforward way and they can handle a mixture of sets and bags.

Let  $N = (P, T, F)$  be a P/T-net. Elements of  $P \cup T$  are called *nodes*. A node  $x$  is an *input node* of another node  $y$  iff there is a directed arc from  $x$  to  $y$  (i.e.,  $(x, y) \in F$  or  $xFy$  for short). Node  $x$  is an *output node* of  $y$  iff  $yFx$ . For any  $x \in P \cup T$ ,  $\bullet x = \{y \mid yFx\}$  and  $x \bullet = \{y \mid xFy\}$ ; the superscript  $N$  may be omitted if clear from the context.

Figure 4 shows a P/T-net consisting of seven places and six transitions. Transition  $A$  has one input place and two output places. Transition  $A$  is an *AND-split*. Transition  $D$  has two input places and one output place. Transition  $D$  is an *AND-join*. The black dot in the input place of  $A$  and  $E$  represents a token. This token denotes the initial marking. The dynamic behavior of such a marked P/T-net is defined by a *firing rule*.

**Definition 2 (Firing rule)** Let  $N = ((P, T, F), s)$  be a marked P/T-net. Transition  $t \in T$  is enabled, denoted  $(N, s)[t]$ , iff  $\bullet t \leq s$ . The firing rule  $-[-]- \subseteq \mathcal{N} \times T \times \mathcal{N}$  is the smallest relation satisfying for any  $(N = (P, T, F), s) \in \mathcal{N}$  and any  $t \in T$ ,  $(N, s)[t] \Rightarrow (N, s - \bullet t + t \bullet)$ .

In the marking shown in Fig. 4 (i.e., one token in the source place), transitions  $A$  and  $E$  are enabled. Although both are enabled only one can fire. If transition  $A$  fires, a token is removed from its input place and tokens are put in its output places. In the resulting marking, two transitions are enabled:  $B$  and  $C$ . Note that  $B$  and  $C$  can be fired concurrently and we assume *interleaving semantics*. In other words, parallel tasks are assumed to be executed in some order.

**Definition 3 (Reachable markings)** Let  $(N, s_0)$  be a marked P/T-net in  $\mathcal{N}$ . A marking  $s$  is *reachable* from the initial marking  $s_0$  iff there exists a sequence of enabled transitions whose firing leads from  $s_0$  to  $s$ . The set of reachable markings of  $(N, s_0)$  is denoted  $[N, s_0]$ .

The marked P/T-net shown in Fig. 4 has six reachable markings. Sometimes it is convenient to know the sequence of transitions that are fired in order to reach some given marking. This paper uses the following notations for sequences. Let  $A$  be some alphabet of identifiers. A *sequence of length  $n$* , for some natural number  $n \in \mathbb{N}$ , over alphabet  $A$  is a function  $\sigma: \{0, \dots, n - 1\} \rightarrow A$ . The sequence of length zero is called the empty sequence and written  $\varepsilon$ . For the sake of readability, a sequence of positive length is usually written by juxtaposing the function values. For example, a sequence  $\sigma = \{(0, a), (1, a), (2, b)\}$ ,

for  $a, b \in A$ , is written  $aab$ . The set of all sequences of arbitrary length over alphabet  $A$  is written  $A^*$ .

**Definition 4** (*Firing sequence*) Let  $(N, s_0)$  with  $N = (P, T, F)$  be a marked P/T net. A sequence  $\sigma \in T^*$  is called a firing sequence of  $(N, s_0)$  if and only if, for some natural number  $n \in \mathbb{N}$ , there exist markings  $s_1, \dots, s_n$  and transitions  $t_1, \dots, t_n \in T$  such that  $\sigma = t_1 \dots t_n$  and, for all  $i$  with  $0 \leq i < n$ ,  $(N, s_i)[t_{i+1}]$  and  $s_{i+1} = s_i - \bullet t_{i+1} + t_{i+1} \bullet$ . (Note that  $n = 0$  implies that  $\sigma = \varepsilon$  and that  $\varepsilon$  is a firing sequence of  $(N, s_0)$ .) Sequence  $\sigma$  is said to be enabled in marking  $s_0$ , denoted  $(N, s_0)[\sigma]$ . Firing the sequence  $\sigma$  results in a marking  $s_n$ , denoted  $(N, s_0)[\sigma](N, s_n)$ .

For the marked Petri net shown in Fig. 4, some possible firing sequences are  $ABCD$ ,  $ACBD$ , and  $AE$ . Note that, for these firing sequences, the resulting marking has a single token and this token is in the output place of transitions  $D$  and  $F$ .

### 3 Internal representation and semantics

When defining the internal representation to be used by our genetic algorithm, the main requirement was that this representation should express the dependencies between the tasks in the log. In other words, the model should clearly express which tasks would enable the execution of other tasks. Additionally, it would be nice if the internal representation would be compatible with a formalism to which analysis techniques and tools exist. This way, these techniques could also be applied to the discovered models. Thus, one option would be to directly represent the individual (or process model) as a Petri net (Desel and Esparza 1995; Murata 1989). However, such a representation would require determining the number of places in every individual and this is not the core concern. It is more important to show the dependencies between the tasks and the semantics of the split/join tasks. Therefore, we defined an internal representation, that is, as expressive as Petri nets (from the task dependency perspective) but that only focuses on the tasks. This representation is called *causal matrix*. Figure 5 shows in (1) the causal matrix that expresses the same task dependencies that are in the “original Petri net.” The causal matrix shows which tasks enable the execution of other tasks via the matching of *input* ( $I$ ) and *output* ( $O$ ) condition functions. The sets returned by the condition functions  $I$  and  $O$  have *subsets* that contain the tasks in the model. Tasks in a same subset have an XOR-split/join relation. Sets in different subsets have an AND-split/join relation. Thus, every  $I$  and  $O$  set expresses a conjunction of exclusive disjunctions. Additionally, a task may appear in more than one subset in a same set. As an example, for task  $D$  in the original Petri net in Fig. 5 the causal matrix states that  $I(D) = \{\{F, B, E\}, \{E, C\}, \{G\}\}$  because  $D$  is enabled by an AND-join construct that has three places. From top to bottom, the first place has a token whenever  $F$  or  $B$  or  $E$  fires. The second place, whenever  $E$  or  $C$  fires. The third place, whenever  $G$  fires. Similarly, the causal matrix has  $O(D) = \{\}$  because  $D$

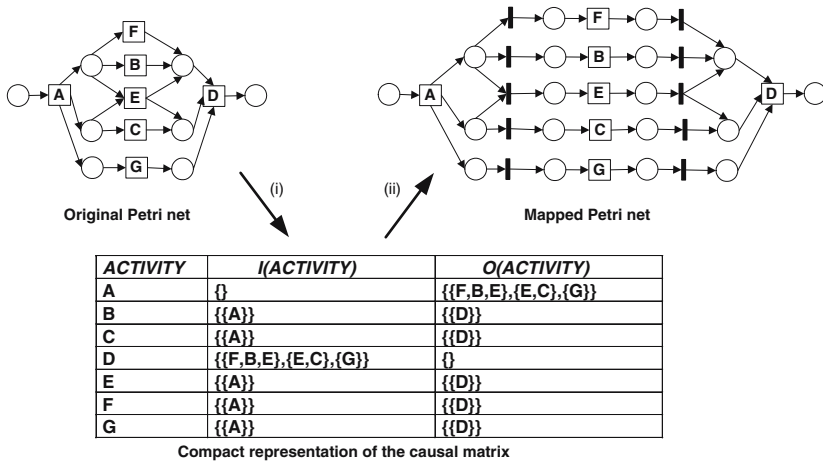


Fig. 5 Mapping of a PN with more than one place between two tasks (or transitions)

is executed last in the model. The following definition formally defines these notions.

**Definition 5 (Causal matrix)** A Causal Matrix is a tuple  $CM = (A, C, I, O)$ , where

- $A$  is a finite set of activities,
- $C \subseteq A \times A$  is the causality relation,
- $I : A \rightarrow \mathcal{P}(\mathcal{P}(A))$  is the input condition function,<sup>3</sup>
- $O : A \rightarrow \mathcal{P}(\mathcal{P}(A))$  is the output condition function,

such that

- $C = \{(a_1, a_2) \in A \times A \mid a_1 \in \bigcup I(a_2)\}$ ,<sup>4</sup>
- $C = \{(a_1, a_2) \in A \times A \mid a_2 \in \bigcup O(a_1)\}$ ,
- $C \cup \{(a_o, a_i) \in A \times A \mid a_o \overset{C}{\bullet} = \emptyset \wedge \overset{C}{\bullet} a_i = \emptyset\}$  is a strongly connected graph.

The set of all causal matrices is denoted by  $\mathcal{CM}$ , and a bag of causal matrices is denoted by  $\mathcal{CM}[]$ .

Any Petri net without duplicate tasks and without more than one place with the same input tasks and the same output tasks can be mapped to a causal matrix. Definition 6 formalizes such a mapping. The main idea is that there is a causal relation  $C$  between any two tasks  $t$  and  $t'$  whenever at least one of the *output* places of  $t$  is an *input* place of  $t'$ . Additionally, the  $I$  and  $O$  condition functions are based on the input and output places of the tasks. This is a natural way of mapping because the input and output places of Petri nets actually reflect the conjunction of disjunctions that these sets express.

<sup>3</sup>  $\mathcal{P}(A)$  denotes the powerset of some set  $A$ .

<sup>4</sup>  $\bigcup I(a_2)$  is the union of the sets in set  $I(a_2)$ .

**Definition 6** ( $\Pi_{PN \rightarrow CM}$ ) Let  $PN = (P, T, F)$  be a Petri net. The mapping of PN is a tuple  $\Pi_{PN \rightarrow CM}(PN) = (A, C, I, O)$ , where

- $A = T$ ,
- $C = \{(t_1, t_2) \in T \times T \mid t_1 \bullet \cap \bullet t_2 \neq \emptyset\}$ ,
- $I \in T \rightarrow \mathcal{P}(\mathcal{P}(T))$  such that  $\forall_{t \in T} I(t) = \{\bullet p \mid p \in \bullet t\}$ ,
- $O \in T \rightarrow \mathcal{P}(\mathcal{P}(T))$  such that  $\forall_{t \in T} O(t) = \{p \bullet \mid p \in t \bullet\}$ .

The semantics of the causal matrix can be easily understood by mapping them back to Petri nets. This mapping is formalized in Definition 7. Conceptually, the causal matrix behaves as a Petri net that contains visible and invisible tasks. For instance, see Fig. 5. This figure shows (1) the mapping of a Petri net to a causal matrix and (2) the mapping from the causal matrix to a Petri net. The firing rule for the mapped Petri net is very similar to the firing rule of Petri nets in general (cf. Definition 2). The only difference concerns the invisible tasks. Enabled invisible tasks can only fire if their firing enables a visible task. Similarly, a visible task is enabled if all of its input places have tokens or if there exists a set of invisible tasks that are enabled and whose firing will lead to the enabling of the visible task. Conceptually, the causal matrix keeps track of the distribution of tokens at a marking in the output places of the visible tasks. The invisible tasks can be seen as “channels” or “pipes” that are only used when a visible task needs to fire. Every causal matrix starts with a token at the start place. Finally, we point out that, in Fig. 5, although the mapped Petri net does not have the same structure of the original Petri net, these two nets are behaviorally equivalent. In other words, given that these two nets initially have a single token and this token is at the start place (i.e., the input place of  $A$ ), the set of traces the two nets can generate is the same.

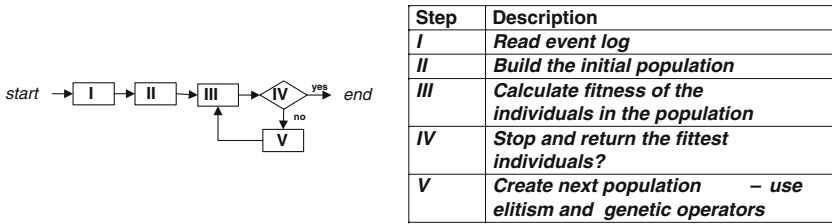
**Definition 7** ( $\Pi_{CM \rightarrow PN}^N$ ) Let  $CM = (A, C, I, O)$  be a causal matrix. The naive Petri net mapping of CM is a tuple  $\Pi_{CM \rightarrow PN}^N = (P, T, F)$ , where

- $P = \{i, o\} \cup \{i_{t,s} \mid t \in A \wedge s \in I(t)\} \cup \{o_{t,s} \mid t \in A \wedge s \in O(t)\}$ ,
- $T = A \cup \{m_{t_1,t_2} \mid (t_1, t_2) \in C\}$ ,
- $F = \{(i, t) \mid t \in A \wedge \bullet t = \emptyset\} \cup \{(t, o) \mid t \in A \wedge t \bullet = \emptyset\} \cup \{(i_{t,s}, t) \mid t \in A \wedge s \in I(t)\} \cup \{(t, o_{t,s}) \mid t \in A \wedge s \in O(t)\} \cup \{(o_{t_1,s}, m_{t_1,t_2}) \mid (t_1, t_2) \in C \wedge s \in O(t_1) \wedge t_2 \in s\} \cup \{(m_{t_1,t_2}, i_{t_2,s}) \mid (t_1, t_2) \in C \wedge s \in I(t_2) \wedge t_1 \in s\}$ .

Definition 7 shows a rather naive approach to generate the mapped Petri net shown in Fig. 5. However, as shown in (Alves de Medeiros et al. 2004b), there are special situations in which more sophisticated mappings are possible.

### 4 Genetic algorithm

In this section, we describe the main steps of our genetic algorithm. Figure 6 shows how they are related.



**Fig. 6** Main steps of our genetic algorithm

### 4.1 Initial population

The initial population is randomly built by the genetic algorithm. As explained in Sect. 3, individuals are causal matrices. When building the initial population, we ensure that the individuals comply with Definition 5. Given a log, all individuals in any population of the genetic algorithm have the same set of activities (or tasks)  $A$ . This set contains the tasks that appear in the log. The setting of the causality relation  $C$  can be done via a completely random approach or a heuristic one. The random approach uses 50% probability for establishing (or not) a causality relation between two task in  $A$ . The heuristic approach uses the information in the log to determine the probability that two tasks are going to have a causality relation set. In a nutshell, the heuristics works as follows: the more often a task  $t_1$  is directly followed by a task  $t_2$  (i.e., the subtrace “ $t_1, t_2$ ” appears in traces in the log), the higher the probability that individuals are built with a causality relation from  $t_1$  to  $t_2$  (i.e.,  $(t_1, t_2) \in C$ ). This heuristic way of building an individual is based on the work presented in Weijters and van der Aalst (2003). Section 4.1.1 has more details about the heuristic approach. Once the causality relations of an individual are determined, the condition functions  $I$  and  $O$  are randomly built. This is done by setting a maximum size  $n$  for any input or output condition function set of a task  $t$  in the initial population.<sup>5</sup> Every task  $t_1$  that causally precedes a task  $t_2$ , i.e.,  $(t_1, t_2) \in C$ , is randomly inserted in *one or more* subsets of the input condition function of  $t_2$ . A similar process is done to set the output condition function of a task.<sup>6</sup> In our case, we set the number of distinct tasks in the log as the maximum size for any input/output condition function set in the initial population.<sup>7</sup> As a result, the initial population can have any individual in the search space defined by a set of activities  $A$ , and that satisfy the constraints for the size of the input/output condition function sets. Note that the higher the amount of tasks that a log contains, the bigger this search space. Finally, we emphasize that no further limitations to the input/output condition functions sets are made in the other steps of the genetic algorithm. Therefore, during the “Step V” in Fig. 6, these sets can increase or shrink as the population evolves.

<sup>5</sup> Formally:  $\forall t \in A [|I(t)| \leq n \wedge |O(t)| \leq n]$ .

<sup>6</sup> Formally:  $\forall t_1, t_2 \in A, (t_1, t_2) \in C [\exists i \in I(t_2) : t_1 \in i]$  and  $\forall t_1, t_2 \in A, (t_1, t_2) \in C [\exists o \in O(t_1) : t_2 \in o]$ .

<sup>7</sup> Formally:  $\forall t \in A [|I(t)| \leq |A| \wedge |O(t)| \leq |A|]$ .

### 4.1.1 Heuristics to build the causality relation of a causal matrix

When applying a genetic algorithm to a domain, it is common practice to “give a hand” to the genetic algorithm by using well-know heuristics (in this domain) to build the initial population (Eiben and Smith 2003). Studies show that the use of heuristics often does not alter the end result (if the genetic algorithm runs for infinite amount of time), but it may speed the early stages of the evolution. The GAs that use heuristics are called *hybrid genetic algorithms*.

In our specific domain—process mining—some heuristics have proven to give reasonable solutions when used to mine event logs. These heuristics are mostly based on local information in the event-log. Due to its similarities to other related work, we use the heuristics in Weijters and van der Aalst (2003) to guide the setting of the causality relations in the individuals of the initial population. These heuristics are based on the *dependency measure*. To define this measure, we first need to formalize the notion of an *event log*.

**Definition 8** (*Event trace, event log*) Let  $T$  be a set of tasks.  $\sigma \in T^*$  is an event trace and  $L: T^* \rightarrow \mathbb{N}$  is an event log. For any  $\sigma \in \text{dom}(L)$ ,  $L(\sigma)$  is the number of occurrences of  $\sigma$ . The set of all event logs is denotes by  $\mathcal{L}$ .

Note that we use  $\text{dom}(f)$  and  $\text{rng}(f)$  to, respectively, denote the *domain* and *range* of a function  $f$ . Furthermore, we use the notation  $\sigma \in L$  to denote  $\sigma \in \text{dom}(L) \wedge L(\sigma) \geq 1$ . For example, assume a log  $L = [abcd, acbd, abcd]$  for the net in Fig. 4. Then, we have that  $L(abcd) = 2$ ,  $L(acbd) = 1$ , and  $L(ab) = 0$ .

The dependency measure basically indicates how strongly a task depends (or is caused) by another task. The more often a task  $t_1$  *directly* precedes another task  $t_2$  in the log, and the less often  $t_2$  *directly* precedes  $t_1$ , the stronger is the dependency between  $t_1$  and  $t_2$ . In other words, the more likely it is that  $t_1$  is a *cause* to  $t_2$ . The dependency measure is given in Definition 9. The notation used in this definition is as follows.  $l2l: T \times T \times \mathcal{L} \rightarrow \mathbb{N}$  is a function that detects length-two loops.  $l2l$  gives the number of times that the substring “ $t_1t_2t_1$ ” occurs in the log  $L$ .  $follows: T \times T \times \mathcal{L} \rightarrow \mathbb{N}$  is a function that returns the number of times that a task is directly followed by another one. That is, how often the substring “ $t_1t_2$ ” occurs in the log  $L$ .

**Definition 9** (*Dependency measure—D*) Let  $L$  be an event log. Let  $T$  be the set of tasks in  $L$ . Let  $t_1$  and  $t_2$  be two tasks in  $T$ . The dependency measure  $D: T \times T \times \mathcal{L} \rightarrow \mathbb{R}$  is a function defined as:

$$D(t_1, t_2, L) = \begin{cases} \frac{l2l(t_1, t_2, L) + l2l(t_2, t_1, L)}{l2l(t_1, t_2, L) + l2l(t_2, t_1, L) + 1}, & \text{if } t_1 \neq t_2 \text{ and } l2l(t_1, t_2, L) > 0, \\ \frac{follows(t_1, t_2, L) - follows(t_2, t_1, L)}{follows(t_1, t_2, L) + follows(t_2, t_1, L) + 1}, & \text{if } t_1 \neq t_2 \text{ and } l2l(t_1, t_2, L) = 0, \\ \frac{follows(t_1, t_2, L)}{follows(t_1, t_2, L) + 1}, & \text{if } t_1 = t_2. \end{cases}$$

Observe that the dependency relation distinguishes between tasks in short loops (length-one and length-two loops) and tasks in parallel. Moreover, the “+1” in the denominator is used to benefit more frequent observations over less frequent ones. For instance, if a length-one-loop “ $t$ ” happens only once in the log  $L$ , the dependency measure  $D(t, t, L) = 0.5$ . However, if this same length-one-loop would occur a 100 times in the log,  $D(t, t, L) = 0.99$ . Thus, the more often a substring (or pattern) happens in the log, the stronger the dependency measure.

Once the dependency relations are set for the input event log, the genetic algorithm uses it to randomly build the causality relations for every individual in the initial population. The pseudo-code for this procedure is the following:

**Pseudo-code:**

**input:** An event-log  $L$ , a power value  $p$ , the dependency function  $D$ .

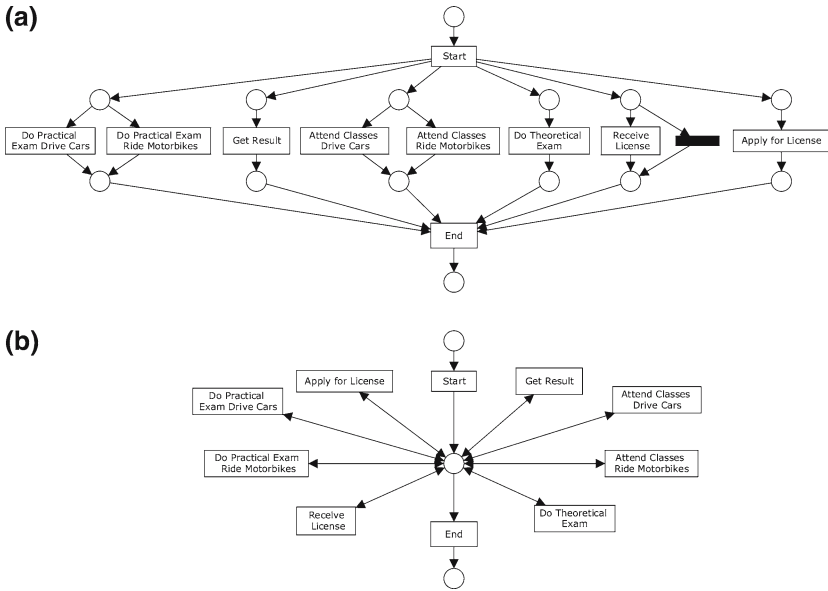
**output:** A causality relation  $C$ .

- (1)  $T \leftarrow$  set of tasks in  $L$ .
- (2)  $C \leftarrow \emptyset$ .
- (3) FOR every tuple  $(t_1, t_2)$  in  $T \times T$  do:
  - (a) Randomly select a number  $r$  between 0 (inclusive) and 1.0 (exclusive).
  - (b) IF  $r < D(t_1, t_2, L)^p$  then:
    - (i)  $C \leftarrow C \cup \{(t_1, t_2)\}$ .
- (4) Return the causality relation  $C$ .

Note that we use a power value  $p$  to control the “influence” of the dependency measure in the probability of setting a causality relation. Higher values for  $p$  lead to the inference of fewer causality relations among the tasks in the event log, and vice-versa.

## 4.2 Fitness calculation

As discussed in Sect. 1, process mining aims at discovering a process model from an event log. This mined process model should give a good insight about what the behavior in the log is. In other words, the mined process model should be *complete and precise from a behavioral perspective*. A process model is complete when it can parse (or reproduce) all the event traces in the log. A process model is precise when it cannot parse more than the traces in the log. The requirement that the mined model should also be precise is important because different models are able to parse all event traces and these models may allow for extra behavior that does not belong to the log. To illustrate this we consider the nets shown in Fig. 7. These models can also parse the traces in Table 1, but they allow for extra behavior. For instance, both models allow for the applicant to take the exam before attending to classes. The fitness function guides the search process of the genetic algorithm. Thus, the fitness of an individual is assessed by benefiting the individuals that can parse more event traces in the log (the “completeness” requirement) and by punishing the individuals that allow for more extra behavior than the one expressed in the log (the “preciseness” requirement).



**Fig. 7** Example of nets that can also reproduce the behavior for the log in Table 1. The problem here is that these nets allow for extra behavior that is not in the log

To facilitate the explanation of our fitness measure, we divide it into three parts. First, we discuss in Sect. 4.2.1, how we defined the part of the fitness measure that guides the genetic algorithm toward individuals that are more complete. Second, we show in Sect. 4.2.2, how we defined the part of the fitness measure that benefits individuals that are more precise. Finally, we show in Sect. 4.2.3 the fitness measure that our genetic algorithm is using. This fitness measure combines the partial fitness measures that are presented in Sects. 4.2.1 and 4.2.2.

### 4.2.1 The “completeness” requirement

The “completeness” requirement of our fitness measure is based on the parsing of event traces by individuals. For a noise-free log, the perfect individual should have fitness 1. This means that this individual could parse all the traces in the log. Therefore, a natural fitness for an individual to a given log seems to be the number of properly parsed event traces<sup>8</sup> divided by the total number of event traces. However, this fitness measure is too coarse because it does not give an indication about (1) how many parts of an individual are correct when the individual does not properly parse an event trace and (2) the semantics of

<sup>8</sup> An event trace is properly parsed by an individual if, for an initial marking that contains a single token and this token is at the start place of the mapped Petri net for this individual, after firing the visible tasks in the order in which they appear in the event trace, the end place is the only one to be marked and it has a single token.



the split/join tasks. For instance, if a net has an AND-split instead of an XOR-split, it may happen that all tasks in a trace can be replayed by this net, but this net does not proper complete for this trace because tokens remain at some of the output places of the AND-split task. So, we defined a more elaborate fitness function: when the task to be parsed is not enabled, the problems (e.g., number of missing tokens to enable this task) are registered and the parsing proceeds as if this task would be enabled. This *continuous* parsing semantics is more robust because it gives a better indication of how many tasks do or do not have problems during the parsing of a trace. The partial fitness function that tackles the “completeness” requirement is in Definition 10. The notation used in this definition is as follows.  $allParsedActivities(L, CM)$  gives the total number of tasks in the event log  $L$  that could be parsed without problems by the causal matrix (or individual)  $CM$ .  $numActivitiesLog(L)$  gives the number of tasks in  $L$ .  $allMissingTokens(L, CM)$  indicates the number of missing tokens in all event traces.  $allExtraTokensLeftBehind(L, CM)$  indicates the number of tokens that were not consumed after the parsing has stopped plus the number of tokens of the end place minus 1 (because of proper completion).  $numTracesLog(L)$  indicates the number of traces in  $L$ .  $numTracesMissingTokens(L, CM)$  and  $numTracesExtraTokensLeftBehind(L, CM)$ , respectively, indicate the number of traces in which tokens were missing and tokens were left behind during the parsing.

**Definition 10** (*Partial fitness* –  $PF_{complete}$ ) Let  $L$  be a non-empty event log. Let  $CM$  be a causal matrix. Then the partial fitness  $PF_{complete}: \mathcal{L} \times \mathcal{CM} \rightarrow (-\infty, 1]$  is a function defined as:

$$PF_{complete}(L, CM) = \frac{allParsedActivities(L, CM) - punishment}{numActivitiesLog(L)}$$

where

$$punishment = \frac{allMissingTokens(L, CM)}{numTracesLog(L) - numTracesMissingTokens(L, CM) + 1} + \frac{allExtraTokensLeftBehind(L, CM)}{numTracesLog(L) - numTracesExtraTokensLeftBehind(L, CM) + 1}.$$

The partial fitness  $PF_{complete}$  gives a more detailed indication about how fit an individual is to a given log. The function  $allMissingTokens$  penalizes (1) nets with XOR-split where it should be an AND-split and (2) nets with an AND-join where it should be an XOR-join. Similarly, the function  $allExtraTokensLeftBehind$  penalizes (1) nets with AND-split where it should be an XOR-split and (2) nets with an XOR-join where it should be an AND-join. Note that we weigh the impact of the  $allMissingTokens$  and  $allExtraTokensLeftBehind$  functions by, respectively, dividing them by the number of event traces minus the number of event traces with missing and left-behind tokens. The main idea is to promote

individuals that correctly parse the more frequent behavior in the log. Additionally, if two individuals have the same *punishment* value, the one that can parse more tasks has a better fitness because its missing and left-behind tokens impact fewer tasks. This may indicate that this individual has more correct *I* and *O* condition functions than incorrect ones. In other words, this individual is a better candidate to produce offsprings for the next population (see Sect. 4.4).

#### 4.2.2 The “preciseness” requirement

The “preciseness” requirement is based on discovering how much extra behavior an individual allows for. To define a fitness measure to punish models that express more than it is in the log is especially difficult because we do not have negative examples to guide our search. Note that the event logs show the allowed (positive) behavior, but they do not express the forbidden (negative) one.

One possible solution to punish an individual that allows for undesirable behavior could be to build the *coverability graph* (Murata 1989) of the mapped Petri net for this individual and check the fraction of event traces this individual can generate that are not in the log. The traces that express different paths of execution for parallelism are not considered as extra behavior. The main idea in this approach is to punish the individual for every extra event trace it generates. Unfortunately, building the coverability graph is not very practical and it is unrealistic to assume that all possible behavior is present in the log.

Because proving that a certain individual is precise is not practical, we use a simpler solution to guide our genetic algorithm toward solutions that have “less extra behavior.” We check, for every marking, the *number of visible tasks that are enabled*. Individuals that allow for extra behavior tend to have more enabled tasks than individuals that do not. For instance, the nets in Fig. 7 have more enabled tasks in most reachable markings than the net in Fig. 2. The main idea in this approach is to benefit individuals that have a smaller amount of enabled tasks during the parsing of the log. This is the measure we use to define our second partial fitness function  $PF_{\text{precise}}$ , that is, presented in Definition 11. The notation used in this definition is as follows. *allEnabledActivities* ( $L, CM$ ) indicates the number of activities that were enabled during the parsing of the log  $L$  by the causal matrix (or individual)  $CM$ . *allEnabledActivities* ( $L, CM[]$ ) apply *allEnabledActivities* ( $L, CM$ ) (see notation for Definition 10) to every element in the bag of causal matrices (or population)  $CM[]$ . The function *max* (*allEnabledActivities* ( $L, CM[]$ )) returns the *maximum* value of the amount of enabled tasks that individuals in the given population ( $CM[]$ ) had while parsing the log ( $L$ ).

**Definition 11** (*Partial fitness*— $PF_{\text{precise}}$ ) Let  $L$  be a non-empty event log. Let  $CM$  be a causal matrix. Let  $CM[]$  be a bag of causal matrices that contains  $CM$ .

The partial fitness  $\text{PF}_{\text{precise}}: \mathcal{L} \times \mathcal{CM} \times \mathcal{CM}[] \rightarrow [0, 1]$  is a function defined as

$$\text{PF}_{\text{precise}}(L, \text{CM}, \text{CM}[]) = \frac{\text{allEnabledActivities}(L, \text{CM})}{\max(\text{allEnabledActivities}(L, \text{CM}[]))}.$$

The partial fitness  $\text{PF}_{\text{precise}}$  gives an indication of how much extra behavior an individual allows for *in comparison to* other individuals in the same population. The smaller the  $\text{PF}_{\text{precise}}$  of an individual is, the better. This way we avoid over-generalizations.

#### 4.2.3 Fitness—combining the “completeness” and “preciseness” requirements

While defining the fitness measure, we decided that the “completeness” requirement should be more relevant than the “preciseness” one. The reason is that we are only interested in precise models that are also complete. The resulting fitness is defined as follows.

**Definition 12** (*Fitness—F*) Let  $L$  be a non-empty event log. Let  $\text{CM}$  be a causal matrix. Let  $\text{CM}[]$  be a bag of causal matrices that contains  $\text{CM}$ . Let  $\text{PF}_{\text{complete}}$  and  $\text{PF}_{\text{precise}}$  be the respective partial fitness functions given in Definitions 10 and 11. Let  $\kappa$  be a real number greater than 0 and smaller or equal to 1 (i.e.,  $\kappa \in (0, 1]$ ). Then the fitness  $F: \mathcal{L} \times \mathcal{CM} \times \mathcal{CM}[] \rightarrow (-\infty, 1)$  is a function defined as

$$F(L, \text{CM}, \text{CM}[]) = \text{PF}_{\text{complete}}(L, \text{CM}) - \kappa * \text{PF}_{\text{precise}}(L, \text{CM}, \text{CM}[]).$$

The fitness  $F$  weighs (by  $\kappa$ ) the punishment for extra behavior. Thus, if a set of individuals can parse all the traces in the log, the one that allows for less extra behavior will have a higher fitness value. For instance, assume a population with the corresponding individual for the net in Fig. 2 and the corresponding individuals for the nets in Fig. 7. If we calculate the fitness  $F$  of these three individuals with respect to the log in Table 1, the individual in Fig. 2 will have the highest fitness value among the three and the individual in Fig. 7(b), the lowest fitness value.

#### 4.3 Stop criteria

The mining algorithm stops when (1) it computes  $n$  generations, where  $n$  is the maximum number of generations that is allowed; or (2) the fittest individual has not changed for  $n/2$  generations in a row.

#### 4.4 Genetic operators

We use *elitism*, *crossover* and *mutation* to build the individuals of the next generation. A percentage of the best individuals (the *elite*) is directly copied

to the next population. The other individuals in the population are generated via crossover and mutation. Two parents produce two offsprings. To select one parent, a *tournament* is played in which five individuals in the population are randomly drawn and the fittest one always wins. The crossover and mutation operator are explained, respectively, in Sects. 4.4.1 and 4.4.2.

#### 4.4.1 Crossover

Crossover is a genetic operator that aims at recombining existing material in the current population. In our case, this material is the current causality relations (cf. Definition 5) in the population. Thus, the crossover operator used by our genetic algorithm should allow for the complete search of the space defined by the existing causality relation in a population. Given a set of causality relations, the search space contains all the individuals that can be created by any combination of a subset of the causality relations in the population. Thus, our crossover operator allows an individual to: lose tasks from the subsets in its *I/O* condition functions (but not necessarily causality relations because a same task may be in more than one subset of an *I/O* condition function), add tasks to the subsets in its *I/O* condition functions (again, not necessarily causality relations), exchange causality relations with other individuals, incorporate causality relations that are in the population but are not in the individual, lose causality relations, decrease the number of subsets in its *I/O* condition functions, and/or increase the number of subsets in its *I/O* condition functions. The *crossover rate* determines the probability that two parents undergo crossover. The crossover point of two parents is a randomly chosen task. The pseudo-code for the crossover operator is as follows:

Pseudo-code:

**input:** Two parents ( $parent_1$  and  $parent_2$ ), crossover rate.

**output:** Two possibly recombined offsprings ( $offspring_1$  and  $offspring_2$ ).

- (1)  $offspring_1 \leftarrow parent_1$  and  $offspring_2 \leftarrow parent_2$ .
- (2) With probability “crossover rate” do:
  - (a) Randomly select a task  $t$  to be the crossover point of the offsprings.
  - (b) Randomly select a swap point  $sp_1$  for  $I_1(t)$ .<sup>9</sup> The swap point goes from position 0 to  $n - 1$ , where  $n$  is the number of subsets in the condition function  $I_1(t)$ .
  - (c) Randomly select a swap point  $sp_2$  for  $I_2(t)$ .
  - (d)  $remainingSet_1(t)$  equals subsets in  $I_1(t)$  that are between position 0 and  $sp_1$  (exclusive).
  - (e)  $swapSet_1(t)$  equals subsets in  $I_1(t)$  whose position equals or bigger than  $sp_1$ .

<sup>9</sup> We use the notation  $I_j(t)$  to get the subset returned by the input condition function  $I$  of task  $t$  in individual  $j$ . In this pseudo-code, the individuals are the offsprings ( $offspring_1$  and  $offspring_2$ ) and  $j \in \{1, 2\}$ .

- (f) Repeat steps 2d and e but, respectively, use  $remainingSet_2(t)$ ,  $I_2(t)$ ,  $sp_2$  and  $swapSet_2(t)$  instead of  $remainingSet_1(t)$ ,  $I_1(t)$ ,  $sp_1$  and  $swapSet_1(t)$ .
  - (g) FOR every subset  $S_2$  in  $swapSet_2(t)$  do:
    - (1) With equal probability perform *one of the following steps*:
      - (A) Add  $S_2$  as a new subset in  $remainingSet_1(t)$ .
      - (B) Join  $S_2$  with an existing subset  $X_1$  in  $remainingSet_1(t)$ .
      - (C) Select a subset  $X_1$  in  $remainingSet_1(t)$ , remove the elements of  $X_1$  that are also in  $S_2$  and add  $S_2$  to  $remainingSet_1(t)$ .
  - (h) Repeat step 2g but, respectively, use  $S_1$ ,  $swapSet_1(t)$ ,  $X_2$  and  $remainingSet_2(t)$  instead of  $S_2$ ,  $swapSet_2(t)$ ,  $X_1$  and  $remainingSet_1(t)$ .
  - (i)  $I_1(t) \leftarrow remainingSet_1(t)$  and  $I_2(t) \leftarrow remainingSet_2(t)$ .
  - (j) Repeat steps 2b–h but use  $O(t)$  instead of the  $I(t)$ .
  - (k) Update the related tasks to  $t$ .
- (3) Return  $offspring_1$  and  $offspring_2$ .

Note that, after crossover, the number of causality relations for the whole population remains constant, but how these relations appear in the offsprings may be different from the parents. Moreover, the offsprings may be different even when both parents are equal. For instance, consider the situation in which the crossover operator receives as input two parents that are equal to the causal matrix in Fig. 5. Assume that (1) the crossover point is the task  $D$ , (2) we are doing crossover over the input condition function  $I(D) = \{\{F, B, E\}, \{E, C\}, \{G\}\}$ , and (3) the swap points are  $sp_1 = 1$  and  $sp_2 = 2$ . Then, we have that the  $remainingSet_1(D) = \{\{F, B, E\}\}$ , the  $swapSet_1(D) = \{\{E, C\}, \{G\}\}$ , the  $remainingSet_2(D) = \{\{F, B, E\}, \{E, C\}\}$ , the  $swapSet_2(D) = \{\{G\}\}$ . Let us first crossover the subsets in the  $swapSet_2(D)$  with the  $remainingSet_1(D)$ . During the crossover, the genetic algorithm randomly chooses to merge the subset  $S_2 = \{G\}$  in the  $swapSet_2(D)$  with the existing subset  $X_1 = \{F, B, E\}$ . In a similar way, while swapping the subsets in  $swapSet_1(D)$  with the  $remainingSet_2(D)$ , the algorithm randomly chooses (1) to insert the subset  $S_1 = \{E, C\}$  and remove task  $E$  from the subset  $X_2 = \{F, B, E\}$ , and (2) to insert the subset  $S_1 = \{G\}$  as a new subset in the  $remainingSet_2(D)$ . The result is that  $I_1(D) = \{\{F, B, E, G\}\}$  and  $I_2(D) = \{\{F, B\}, \{E, C\}, \{G\}\}$ . The output condition functions  $O_1(D)$  and  $O_2(D)$  do not change after the crossover operator because the task  $D$  does not have any output task. After the crossover, the mutation operator takes place.

#### 4.4.2 Mutation

The mutation operator aims at inserting new material in the current population. In our case, this means that the mutation operator may change the existing causality relations of a population. Thus, our mutation operator performs one of the following actions to the  $I/O$  condition functions of a task in an individual: (1) randomly choose a subset and add a task (in  $A$ ) to this subset, (2) randomly choose a subset and remove a task out of this subset, or (3) randomly redistribute the elements in the subsets of  $I/O$  into new subsets. For example, consider

the input condition function of task  $D$  in Fig. 5.  $I(D) = \{\{F, B, E\}, \{E, C\}, \{G\}\}$  can be mutated to (1)  $\{\{F, B, E\}, \{E, C\}, \{G, D\}\}$  if task  $D$  is added to the subset  $\{G\}$ , (2)  $\{\{F, B, E\}, \{C\}, \{G\}\}$  if task  $E$  is removed from the subset  $\{E, C\}$ , or (3)  $\{\{F\}, \{E, C, B\}, \{G\}, \{E\}\}$  if the elements in the original  $I(D)$  are randomly redistributed in a randomly chosen number of new subsets. Every task in an offspring may undergo mutation with the probability determined by the *mutation rate*. The pseudo-code for the mutation operator is as follows:

**Pseudo-code:**

**input:** An individual, mutation rate.

**output:** A possibly mutated individual.

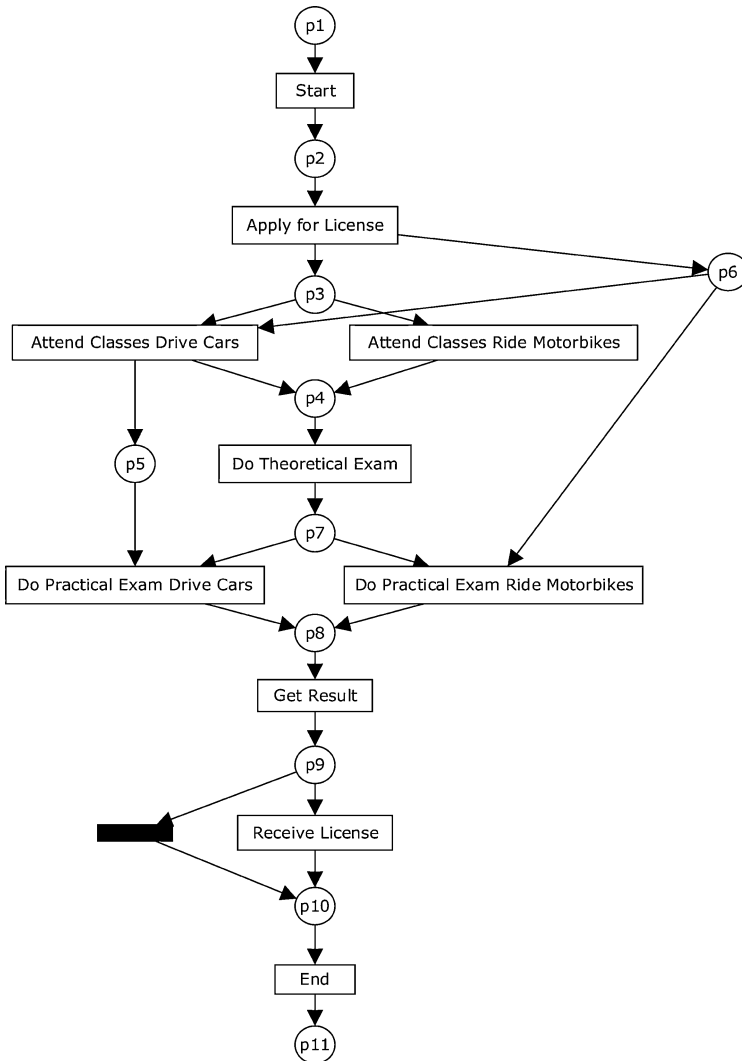
- (1) For every task  $t$  in the individual do:
  - (a) With probability *mutation rate* do one of the following operations for the condition function  $I(t)$ :
    - (1) Select a subset  $X$  in  $I(t)$  and add a task  $t'$  to  $X$ , where  $t'$  belongs to the set of tasks in the individual.
    - (2) Select a subset  $X$  in  $I(t)$  and remove a task  $t'$  from  $X$ , where  $t'$  belongs to  $X$ . If  $X$  is empty after  $t'$  removal, exclude  $X$  from  $I(t)$ .
    - (3) Redistribute the elements in  $I(t)$ .<sup>10</sup>
  - (b) Repeat step 1a, but use the condition function  $O(t)$  instead of  $I(t)$ .
  - (c) Update the related tasks to  $t$ .

As the reader may already have noticed, both the crossover and the mutation operators perform a repairing operation at the end of their executions. The “update the related tasks” operation makes sure that the individual is still compliant with Definition 5 after undergoing crossover and/or mutation.

## 5 Analysis metrics

The genetic algorithm searches for models that are *complete* and *precise* (see Sect. 4.2). Therefore, when evaluating the results of our experiments, we should check if the mined models are indeed complete and precise. At first sight, the natural way to check for this seemed to be to compare the causal matrix of the original model (the one that was simulated to created the synthetic event logs) with the causal matrix of the individual that was mined by the genetic algorithm. However, this is not a good evaluation criterion because there are different ways to model the exact behavior expressed in a log. For instance, consider the net in Fig. 8. This net produces exactly the same behavior as the one in Fig. 2. However, their causal matrices are different. Furthermore, even when the mined models are not complete and/or precise, we should be able to assess how much correct material they contain. This is important because we do not let the genetic algorithm run for an “infinite” amount of time. Thus, even

<sup>10</sup> Details about this step: (1) Get a list with the elements of  $I(t)$ ; (2) Create  $n$  sets ( $n$  is the number of elements in  $I(t)$ ) and randomly distribute them in the  $n$  sets; and (3) Filter out the non-empty sets. These non-empty sets are now the subsets of  $I(t)$ .



**Fig. 8** Other mined net for the log in Table 1. Note that this net is behaviorally equivalent to the net in Fig. 2, although they are structurally different. Note that the place “p6” has different input and output tasks in the two nets

when the mined model is not complete and precise, it is important to know if the genetic algorithm is going in the right direction.

In our experiments, we have three elements: (1) the *original model* that is used to build the synthetic event log, (2) the synthetic *event log* itself, and (3) the *mined model* (or individual). Thus, to analyze our results, we have defined metrics that are based on two or more of these elements.

## 5.1 Checking for completeness

To check for completeness, we only need the event log and the mined model. Recall that a model is complete when it can parse all the traces in the log without having missing tokens or tokens left behind. So, completeness can be verified by calculating the partial fitness  $PF_{\text{complete}}$  (see Definition 10) for the event log and the mined model. Whenever  $PF_{\text{complete}} = 1$ , the mined model is complete. Moreover, even when the mined model has  $PF_{\text{complete}} < 1$ , this measure gives an indication of the quality of the mined model with respect to completeness.

## 5.2 Checking for preciseness

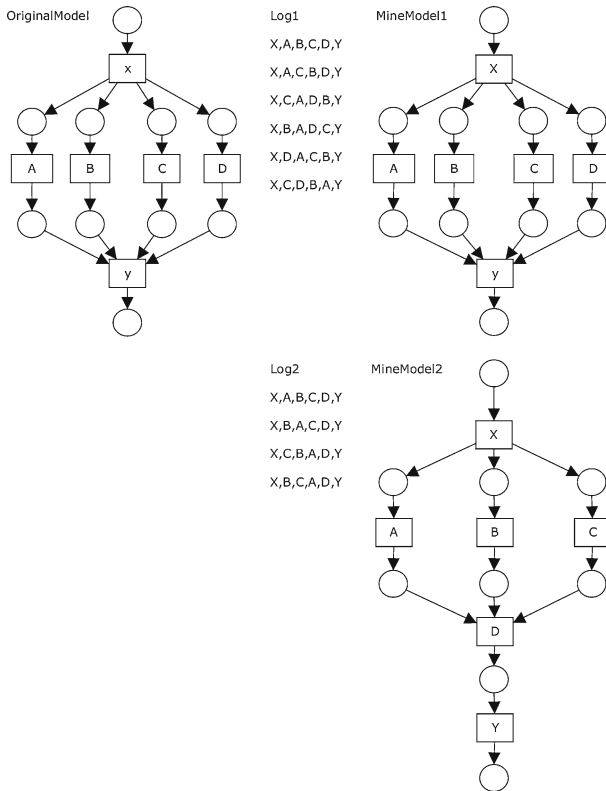
To check for preciseness, we need the original model, the event log and the mined model.<sup>11</sup> The main reason why we could not define metrics only based on the event log and the mined model, or the original model and the mined model, is because it is unrealistic to assume that the event log has all the possible traces that the original model can generate. In other words, it is unrealistic to assume that the log contains all possible event traces. Recall that a model is precise when it does not allow for more behavior than the one expressed in the log. Thus, if the log would be exhaustive, a possible metric to check for this preciseness could be to divide the number of traces that are in the log and that the mined model can generate by the amount of traces that the mined model can generate. Clearly, a precise mined model could not generate more traces than the ones in the log. Note that this metric would be based on the event log and the mined model. Furthermore, metrics based on the mined and original models would also be possible if the log would be entire. For instance, we could compare the coverability graphs (Murata 1989) of mapped Petri nets of the mined and the original models. In this case, the mined model would be precise whenever the coverability graphs would be equal. Note that sophisticated notions such as bisimulation (Milner et al. 1992) and branching bisimulation (van Glabbeek and Weijland 1996) could also be used. However, none of these metrics are suitable because in real-life applications the log does not hold all possible traces.

For instance, consider the situation shown in Fig. 9. This figure shows the original model (“OriginalModel”), two synthetic logs (“Log1” and “Log2”) and their respective mined models (“MinedModel1” and “MinedModel2”). “Log1” shows that the tasks  $A, B, C$ , and  $D$  are (1) always executed after the task  $X$  and before the task  $Y$  and (2) independent of each other. Thus, we can say that the “MinedModel1” is precise with respect to the behavior observed in the “Log1.” However, note that the “MinedModel1,” although precise, can generate more traces than the ones in the “Log1.” A similar reasoning can

---

<sup>11</sup> Note that in reality we do not know the original (or initial) model. However, the only way to evaluate our results is to assume an initial model. Without an initial model, it is impossible to judge preciseness. In other words, there could be over-fitting or over-generalization, but it would be impossible to judge this.





**Fig. 9** Example of two mined models that are complete and precise with respect to the logs, but both mined models can generate more traces than the ones in the log. Additionally, the coverability graph of the “MinedModel2” is different from the one of the “OriginalModel”

be done for the “Log2” and the “MinedModel2.” Moreover, the coverability graph of the “MinedModel2” is different from the one of the “OriginalModel.” Actually, based on “Log2,” the “MinedModel2” is more precise than the “OriginalModel.” This illustrates that, when assessing how close the behavior of the mined and original models are, we have to consider the event log that was used by the genetic algorithm. Therefore, we have defined two metrics to quantify how similar the behavior of the original model and the mined model are *based on the event log used during the mining process*.

The two metrics are the *behavioral precision* ( $B_P$ ) and the *behavioral recall* ( $B_R$ ). Both metrics are based on the parsing of an event log by the mined model and by the original model. The  $B_P$  and  $B_R$  metrics are, respectively, formalized in Definitions 13 and 14. These metrics basically work by checking, for the continuous semantics parsing of every task in every process instance of the event log, how many tasks are enabled in the mined model and how many are enabled in the original model. The more enabled tasks the models have in common, the more similar their behaviors are with respect to the event log. The behavioral

precision  $B_P$  checks how much behavior is allowed by the mined model that is not by the original model. The behavioral recall  $B_R$  checks for the opposite. Additionally, both metrics take into account how often a trace occurs in the log. This is especially important when dealing with logs in which some paths are more likely than others, because deviations corresponding the infrequent paths are less important than deviations corresponding to frequent behavior. Note that, assuming a log generated from an original model and a mined model for this log, we can say that the closer their  $B_P$  and  $B_R$  are to 1, the more similar their behaviors. More specifically, we can say that:

- The mined model is *as precise as* the original model whenever  $B_P$  and  $B_R$  are equal to 1. This is exactly the situation shown in Fig. 9 for the “OriginalModel,” the “Log1” and the “MinedModel1.”
- The mined model is *more precise than* the original model whenever  $B_P = 1$  and  $B_R < 1$ . For instance, see the situation shown in Fig. 9 for the “OriginalModel,” the “Log2” and the “MinedModel2.”
- The mined model is *less precise than* the original model whenever  $B_P < 1$  and  $B_R = 1$ . For instance, see the situation illustrated for the original model in Fig. 2, the log in Fig. 1, and the mined models in Fig. 7.

**Definition 13** (*Behavioral precision— $B_P$* )<sup>12</sup> Let  $L$  be an event log. Let  $CM_o$  and  $CM_m$  be the respective causal matrices for the original (or base) model and for the mined one. Then the behavioral precision  $B_P: \mathcal{L} \times \mathcal{CM} \times \mathcal{CM} \rightarrow [0, 1]$  is a function defined as:

$$B_P(L, CM_o, CM_m) = \frac{\sum_{\sigma \in L} \left( \frac{|L(\sigma)|}{|\sigma|} \times \sum_{i=1}^{|\sigma|} \frac{|Enabled(CM_o, \sigma, i) \cap Enabled(CM_m, \sigma, i)|}{|Enabled(CM_m, \sigma, i)|} \right)}{\sum_{\sigma \in L} |L(\sigma)|}$$

where  $Enabled(CM, \sigma, i)$  gives the enabled activities at the causal matrix  $CM$  just before the parsing of the element at position  $i$  in the trace  $\sigma$ . During the parsing a continuous semantics is used (see Sect. 4.2.1).

**Definition 14** (*Behavioral recall— $B_R$* ) Let  $L$  be an event log. Let  $CM_o$  and  $CM_m$  be the respective causal matrices for the original (or base) model and for the mined one. Then the behavioral recall  $B_R: \mathcal{L} \times \mathcal{CM} \times \mathcal{CM} \rightarrow [0, 1]$  is a

<sup>12</sup> For both Definitions 13 and 14, whenever the denominator “ $|Enabled(CM, \sigma, i)|$ ” is equal to 0, the whole division is equal to 0. For simplicity reasons, we have omitted this condition from the formulae.

function defined as:

$$B_R(L, CM_o, CM_m) = \frac{\sum_{\sigma \in L} \left( \frac{L(\sigma)}{|\sigma|} \times \sum_{i=1}^{|\sigma|} \frac{|Enabled(CM_o, \sigma, i) \cap Enabled(CM_m, \sigma, i)|}{|Enabled(CM_o, \sigma, i)|} \right)}{\sum_{\sigma \in L} L(\sigma)}$$

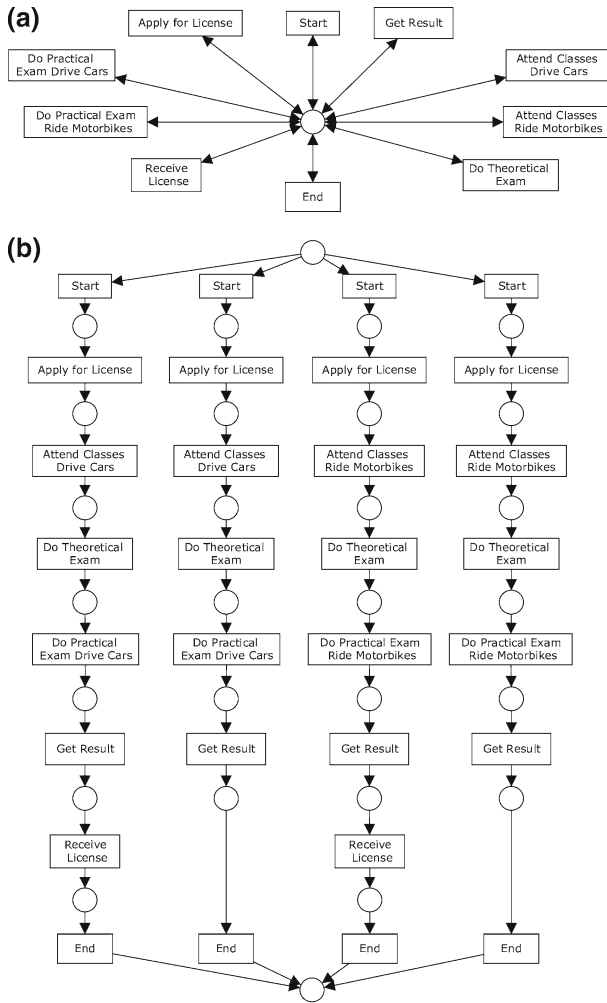
### 5.3 Reasoning about the quality of the mined models

When evaluating the quality of a data mining approach (genetic or not), it is common to check if the approach tends to find over-general or over-specific solutions. In our case, the over-general solution is the one that can parse any trace that can be formed from the tasks in a log. This solution has a self-loop for every task in the log. The over-specific solution is the one that has a branch for every *unique* trace in the log. Figure 10 shows an over-general and an over-specific solution for the log in Table 1.

The over-general solution *does* belong to the search space considered in this paper. However, this kind of solution can be easily detected by the metrics we have defined so far. Note that, for a given original model  $CM_o$ , a log  $L$  generated by simulating  $CM_o$ , and the mined over-general model  $CM_m$ , it always holds that: (1) the over-general model is complete (i.e.,  $PF_{complete}(L, CM_m) = 1$ ); (2) while parsing the traces, all the tasks that are enabled in the original model are also enabled in the over-general model (i.e.,  $B_R(L, CM_o, CM_m) = 1$ ); and (3) while parsing the traces, all the tasks of the over-general model are always enabled, i.e., the formula of the behavioral precision (see Definition 13) can be simplified to the formula in Eq. 1. These three remarks are used to detect over-general mined models during the experiments analysis.

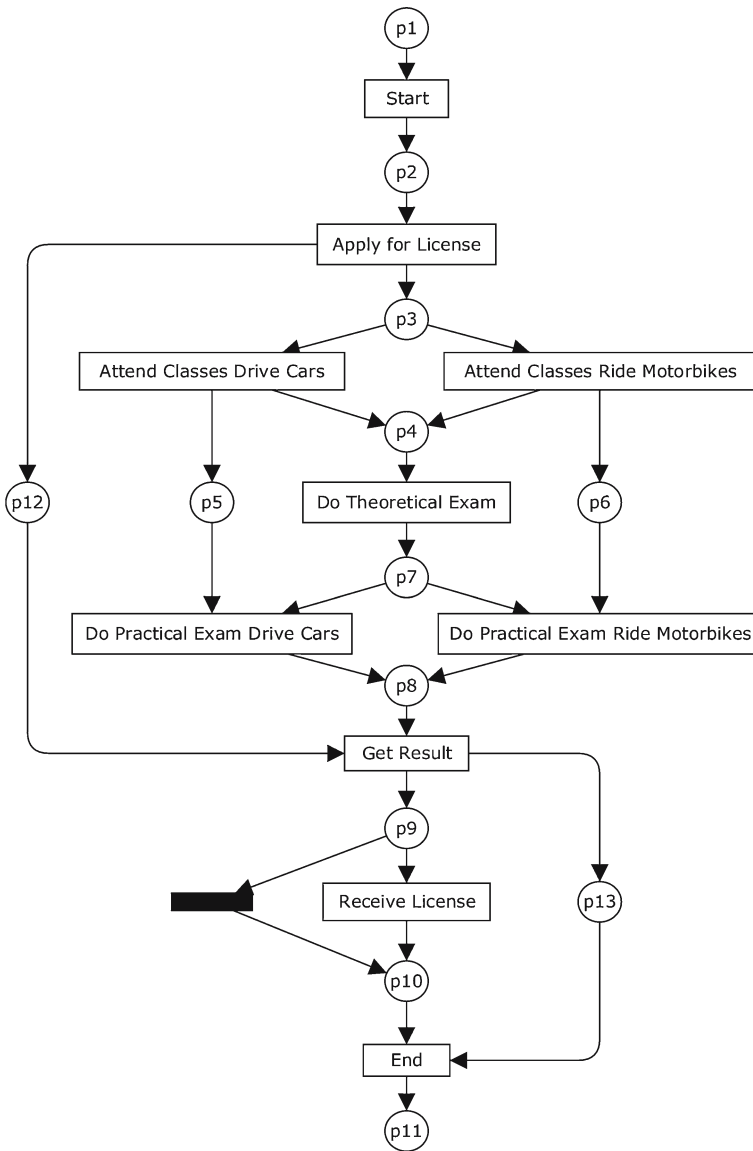
$$B_P(L, CM_o, CM_m) = \frac{\sum_{\sigma \in L} \left( \frac{L(\sigma)}{|\sigma|} \times \sum_{i=1}^{|\sigma|} \frac{|Enabled(CM_o, \sigma, i)|}{|A_m|} \right)}{\sum_{\sigma \in L} L(\sigma)}. \quad (1)$$

Contrary to the over-general solution, the over-specific one *does not* belong to our search space because our internal representation (the causal matrix) does not support duplicate tasks. However, because our fitness only looks for the complete and precise behavior [not the minimal representation, like the works on Minimal Description Length (MDL), Grunwald et al. 2005], it is still important to check how similar the structures of the mined model and the original one are. Differences in the structure may point out another good solution or an overly complex solution. For instance, have a look at the model in Fig. 11. This



**Fig. 10** Example of nets that are (a) over-general and (b) over-specific for the log in the Table 1

net is complete and precise from a behavioral point of view, but it contains extra unnecessary places. Note that the places “p12” and “p13” could be removed from the net without changing its behavior. In other words, “p12” and “p13” are implicit places *van der Aalst et al. (2004)*. Actually, because the places do not affect the net behavior, all the nets in Figs. 2, 8, and 11 have the same fitness. However, a metric that checks the structure of a net would, for instance, point out that the net in Fig. 11 is a “superstructure” of the net in Fig. 2, and has many elements in common with the net in Fig. 8. So, even when we know that the over-specific solution is out of the search space defined in this paper, it is interesting to get a feeling about the structure of the mined models. That is why



**Fig. 11** Example of a net that is behavioral precise and complete w.r.t. the log in Table 1, but that contains extra unnecessary (implicit) places (p12 and p13)

we developed two metrics to assess how much the mined and original model have in common *from a structural point of view*.

The two metrics are the *structural precision* ( $S_P$ ) and the *structural recall* ( $S_R$ ). Both metrics are based on the *causality relations* of the mined and original models, and were adapted from the precision and recall metrics presented in [Pinter and Golani \(2004\)](#). The  $S_P$  and  $S_R$  metrics are, respectively, formalized in

Definitions 15 and 16. These metrics basically work by checking how many causality relations the mined and the original models have in common. The more causality relations the two models have in common, the more similar their structures are. The structural precision assess how many causality relations the mined model has that are not in the original model. The structural recall works the other way around. Note that the structural similarity performed by these metrics does not consider the semantics of the split/join points. We have done so because the causality relations are the core of our genetic material (see Sect. 4.4). The semantics of the split/join tasks can only be correctly captured if the right dependencies (or causality relations) between the tasks in the log are also in place.

**Definition 15** (*Structural precision*— $S_P$ )<sup>13</sup> Let  $CM_o$  and  $CM_m$  be the respective causal matrices for the original and the mined models. The structural precision  $S_P: \mathcal{CM} \times \mathcal{CM} \rightarrow [0, 1]$  is a function defined as:

$$S_P(CM_o, CM_m) = \frac{|C_o \cap C_m|}{|C_m|}.$$

**Definition 16** (*Structural recall*— $S_R$ ) Let  $CM_o$  and  $CM_m$  be the respective causal matrices for the original and the mined model. The structural recall  $S_R: \mathcal{CM} \times \mathcal{CM} \rightarrow [0, 1]$  is a function defined as:

$$S_R(CM_o, CM_m) = \frac{|C_o \cap C_m|}{|C_o|}.$$

When the original and mined models have behavioral metrics  $B_R$  and  $B_P$  that are equal to 1, the  $S_R$  and  $S_P$  show how similar the structure of these models are. For instance, for the original model in Fig. 2, the structural metrics would indicate that a mined model like the one in Fig. 8 differs from the original one by the same amount of causality relations ( $S_R = S_P$ ), and a mined model like the one in Fig. 11 has extra causality relations ( $S_R = 1$  and  $S_P < 1$ ).

#### 5.4 Recapitulation of the analysis metrics

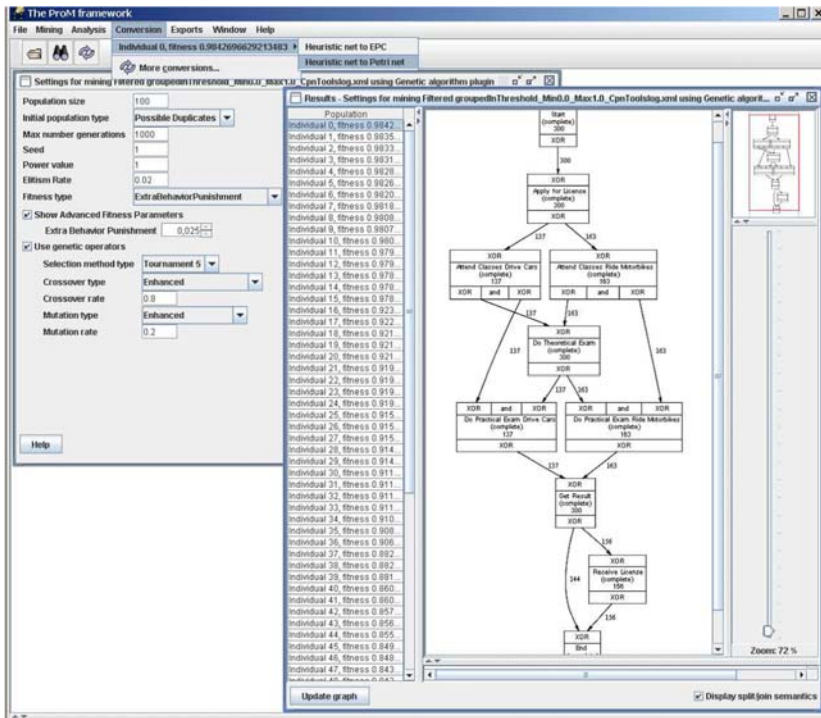
This section has presented the five metrics that are used to analyze the experiments in this paper: the partial fitness for the completeness requirement ( $PF_{\text{complete}}$ ), the behavioral precision ( $B_P$ ), the behavioral recall ( $B_R$ ), the structural precision ( $S_P$ ), and the structural recall ( $S_R$ ). The  $PF_{\text{complete}}$  quantifies how complete a mined model is. The  $B_P$  and  $B_R$  measure how precise the mined model is. The  $S_P$  and  $S_R$  express if the mined model has an overly complex structure or not. These metrics are complementary and should be considered together during the experiments analysis. For instance, for our experiments,

<sup>13</sup> For both Definitions 15 and 16, whenever the denominator “ $|C|$ ” is equal to 0, the whole division is equal to 0. For simplicity reasons, we have omitted this condition from the formulae.

the mined model is as complete and precise as the original model whenever the metrics  $PF_{\text{complete}}$ ,  $B_P$ , and  $B_R$  are equal to 1. More specifically, the mined model is exactly like the original model when all the five metrics are equal to 1. As a general rule, the closer the values of the five metrics are to 1, the better.

## 6 Experiments and results

This section explains how we conducted the experiments and, most important of all, how we analyzed the quality of the models that the genetic algorithm mined. To conduct the experiments we needed (1) to implement our genetic algorithm and (2) a set of event logs. The *genetic algorithm* described in this paper is implemented as the “Genetic algorithm plug-in” in the ProM framework (see Fig. 12). The ProM framework is available at [www.processmining.org](http://www.processmining.org) and supports the development of plug-ins to mine event logs. Although in this paper we focus on the “Genetic algorithm plug-in,” the ProM framework offers other plug-ins like, for instance, the “Social network miner plug-in” (van der Aalst and Song



**Fig. 12** Screenshot of the “Genetic algorithm plug-in” in the ProM framework. This screenshot shows the result of mining an event log like the one in Table 1. This log has 300 process instances in total. On the left-side window shows the configuration parameters (see Sect. 6.1). On the right-side window shows the best mined individual (or causal matrix). Additionally, in the menu bar we show how to convert this individual (called “Heuristics Net” in the ProM framework) to a Petri net

2004) and “Conformance checker plug-in” (Rozinat and van der Aalst 2005). The logs used in our experiments are *synthetic*. In brief, we built the model (or copied it from some related work) and simulated it to create a synthetic event log. We then run the genetic algorithm over these sets of logs. Once the genetic algorithm finished the mining process, we analyzed the results.

When conducting the experiments we were interested in getting an indication of two main points: (1) how the heuristics and genetic operators influence the quality of mined models, and (2) if the fitness measure indeed guides the search toward individuals that reflect the most frequent behavior in the log. The first point was investigated by executing the genetic algorithm over synthetic *noise-free* log. Details about these experiments are explained in Sect. 6.1. The second point was investigated by running the genetic algorithm over synthetic *noisy* logs. Details about the experiments with noisy logs are reported in Sect. 6.2.

## 6.1 Experiments with noise-free logs

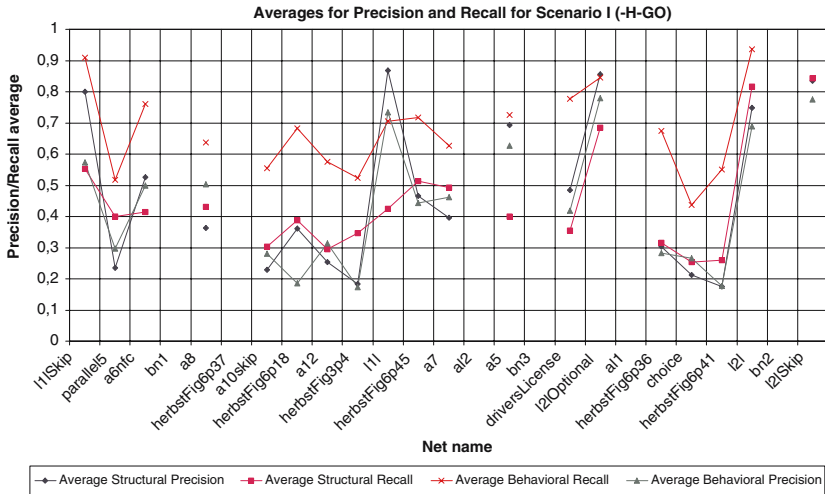
### 6.1.1 Setup

The genetic algorithm was tested over noise-free event logs from 25 different process models. These models contain constructs like sequence, choice, parallelism, loops, non-free-choice, and invisible tasks. From the 25 models, six were copied from the models in (Herbst 2001). The other models were created by the authors. The models had between 6 and 42 tasks.<sup>14</sup> Every event log was randomly generated and contained 300 process instances. To speed up the computational time of the genetic algorithm, the similar traces were grouped into a single one and a counter was associated to inform how often the trace occurs. The similarity criterion was the local information in the trace. Traces with the same direct left and right neighbors for every element were grouped together. Besides, *to test how strong the use of the genetic operators and the heuristics influence the results*, we set up four scenarios while running the genetic algorithm: (“Scenario I”) without heuristics to build the initial population and without genetic operators;<sup>15</sup> (“Scenario II”) with heuristics, but without the genetic operators; (“Scenario III”) without heuristics, but with genetic operators; and (“Scenario IV”) with heuristics and genetic operators. For every log, 50 runs were executed for every scenario. Every run had a population size of 100 individuals, at most 1,000 generations, an elite of two individuals and a  $\kappa$  of 0.025 (see Definition 12). The experiments with heuristics used a power value of 1 while building the initial population (see Sect. 4.1.1). The experiments with the genetic operators have a respective crossover and mutation probabilities of 0.8

<sup>14</sup> Table 5 (cf. Appendix A) provides more details about the characteristics of these models.

<sup>15</sup> This scenario is a random generation of individuals. The aim of experimenting with this scenario is to assess if the use of genetic operators and/or heuristics is better than the pure random generation of individuals, given the same limited amount of computational time.





**Fig. 13** Average precision and recall values of the results for Scenario I (without heuristics to build the initial population and without using genetic operators to build the following populations)

and 0.2 (see the respective Sects. 4.4.1 and 4.4.2). All the experiments were run using the ProM framework, our tool set that can be obtained via [www.process-mining.org](http://www.process-mining.org). We implemented the genetic algorithm and the metrics (cf. Sect. 5) described in this paper as plug-ins for this framework.

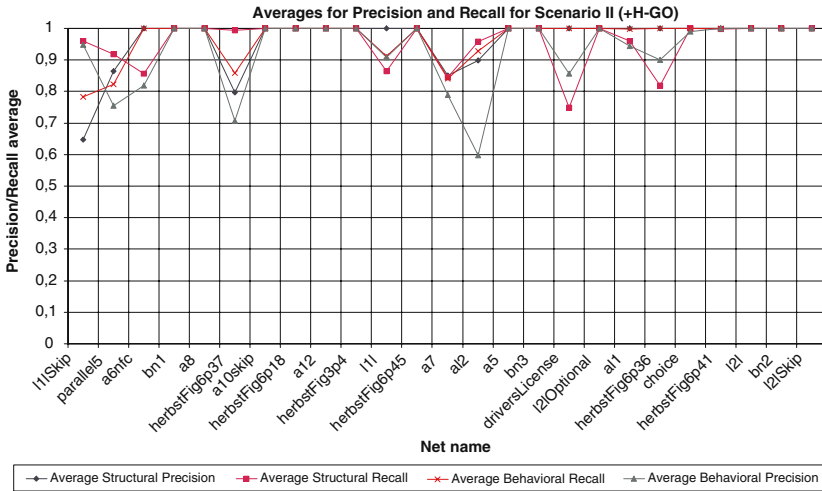
### 6.1.2 Results

The results are reported in Figs. 13–23. Figures 13–19 show the average values of the analysis metrics (cf. Sect. 5) for the mined models. Figures 20 and 21 indicate how the fitness of the best mined model evolved over generations. Figures 22 and 23 show how much computational time was required, on average, per run of the genetic algorithm.<sup>17</sup> Overall, the results indicate that the scenario for the hybrid genetic algorithm (Scenario IV) is superior to the other scenarios in all aspects. More specifically, the results show that:

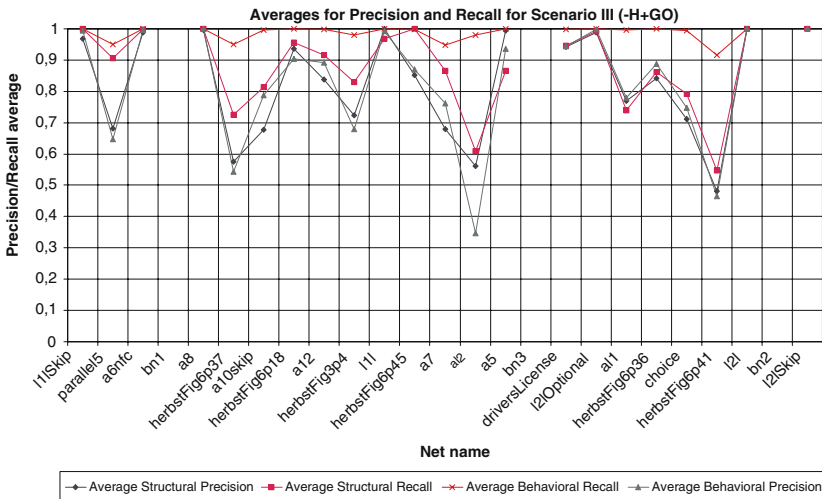
- Any approach (scenarios II–IV) is better than the pure random generation of individuals (Scenario I) (cf. Figs. 13–19).
- Scenarios II and IV mined more complete and precise models than the other scenarios (cf. Figure 19).
- The hybrid genetic algorithm (Scenario IV) works the best. This approach combines the strong ability of the heuristics to correctly capture the *local* causality relations with the benefits of using the genetic operators (especially

<sup>16</sup> The unmarked points in these figures correspond to experiments that were interrupted because they were taking more than 6 hours to process one seed.

<sup>17</sup> The experiments were run in a Intel® Pentium® 4 CPU 3.40 GHz 3.39 GHz, 1.99 GB RAM, with Microsoft Windows XP 2002.

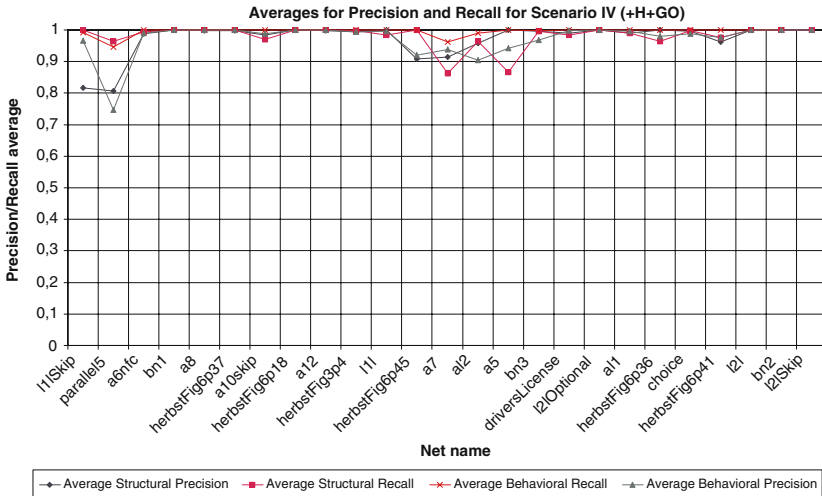


**Fig. 14** Average precision and recall values of the results for Scenario II (with heuristics to build the initial population, but without using genetic operators to build the following populations)

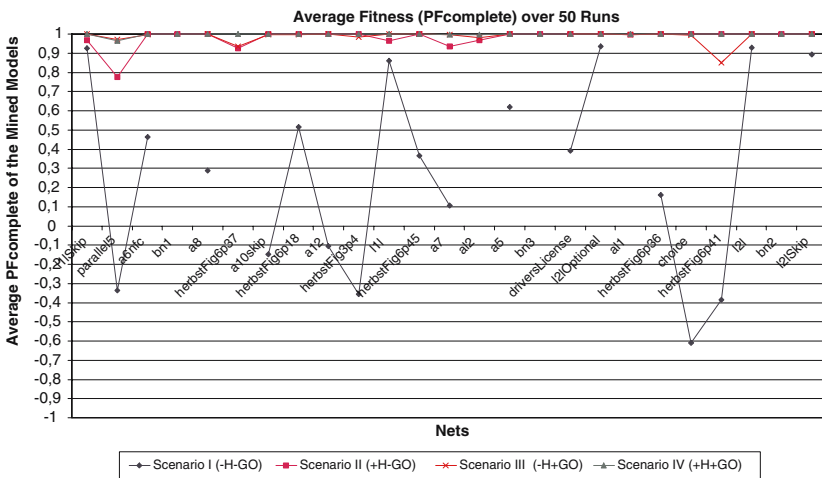


**Fig. 15** Average precision and recall values of the results for Scenario III (without heuristics to build the initial population, but using genetic operators to build the following populations)

mutation) to introduce the non-local causality relations. For instance, consider the results for the nets a6nfc, driversLicense and herbstFig6p36. All these nets have non-local non-free-choice constructs. Note that, for these three nets, the results for Scenario II (cf. Fig. 14) have a much lower behavioral precision than for Scenario IV (cf. Fig. 16). This illustrates the importance of the genetic operators to insert the non-local causality relations.

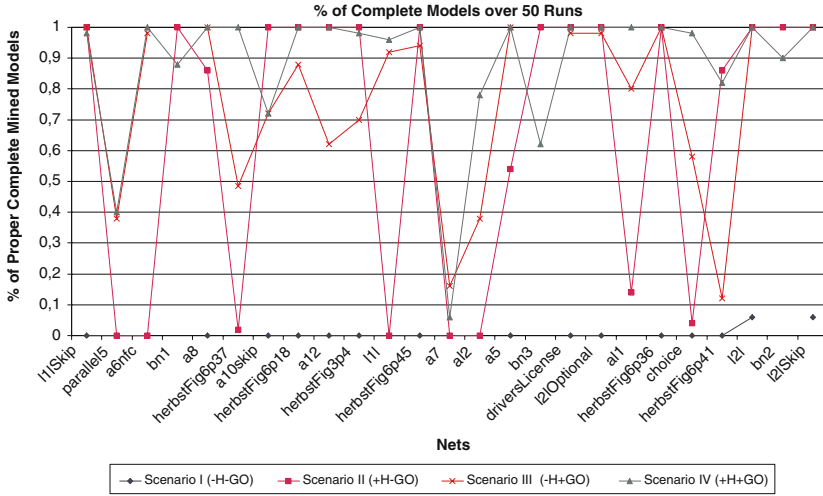


**Fig. 16** Average precision and recall values of the results for Scenario IV (with heuristics to build the initial population and using genetic operators to build the following populations)

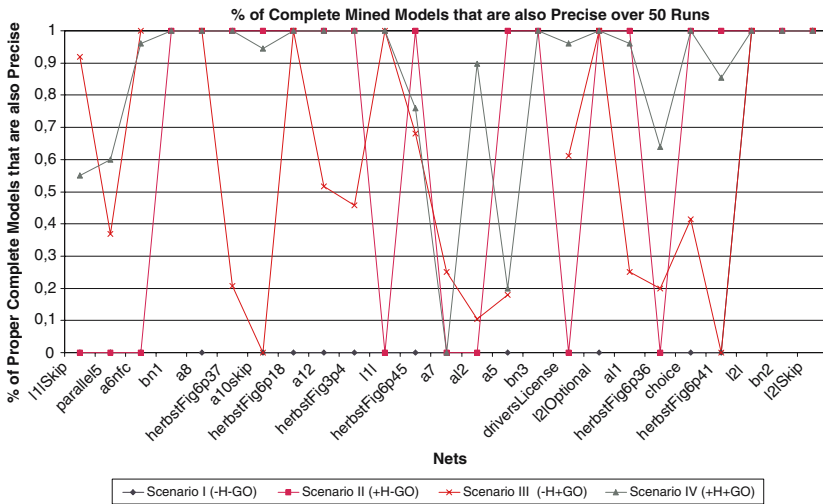


**Fig. 17** Average fitness ( $PF_{complete}$ ) values of the mined models for 50 runs

- In general, the hybrid genetic algorithm (Scenario IV) mines more complete models *that are also precise* than the other scenarios (see Fig. 19). In fact, except for net a7, Scenario IV is the only configuration that mined a complete and precise model for at least one of the runs (cf. Fig. 19). This shows that this scenario finds mined models that are complete and precise faster than the other scenarios.
- Nets with short parallel branches (like parallel15, a7 and a5) are more difficult to mine. This is due to the probabilistic nature of the genetic algorithm. Recall that the fitness measure always benefits the individuals that



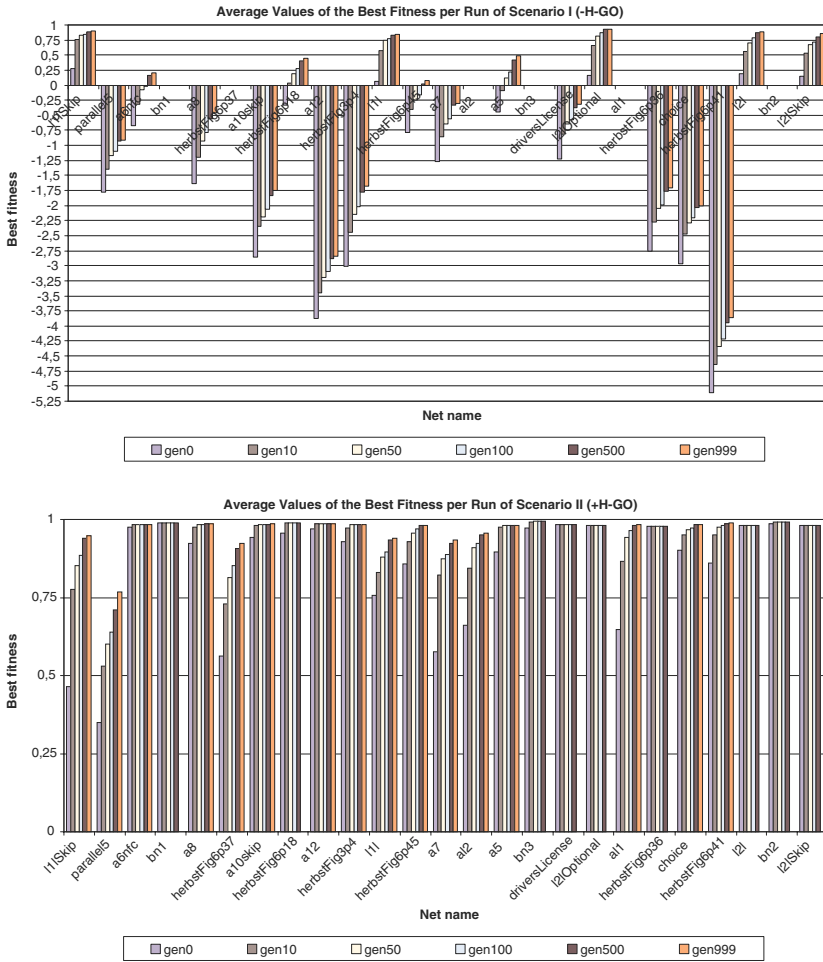
**Fig. 18** Percentage of the mined models that proper complete ( $PF_{complete} = 1$ ) over 50 runs



**Fig. 19** Percentage of the complete mined models ( $PF_{complete} = 1$ ) that are also precise ( $B_p = 1$ ) over 50 runs

can parse the most frequent behavior in the log. So, in parallel situations, it is often the case that the algorithm goes for individuals that show the most frequent interleaving patterns in the log.

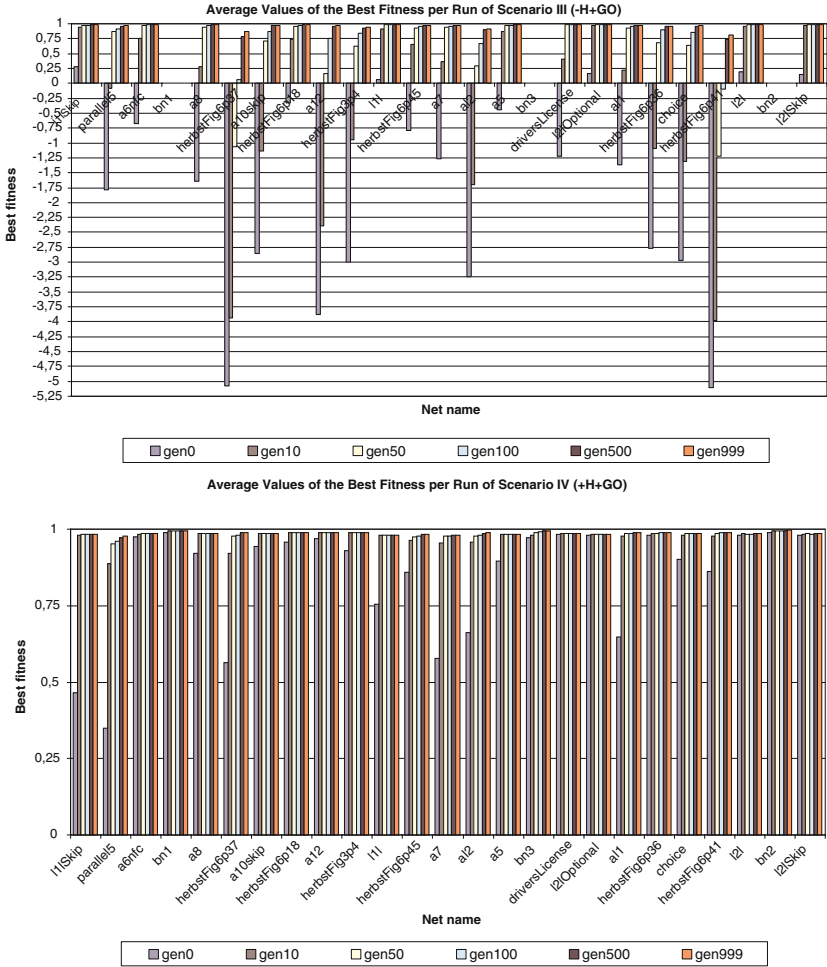
- Although Scenario II led to better results than Scenario III, it is not fair to compare them. The reason is that Scenario III *starts from scratch*, in the sense that its initial population is randomly built, while Scenario II *is strongly helped by good heuristics* to detect local dependencies. In our experiments, Scenario III is used to show that (1) the use of the genetic operators improves the



**Fig. 20** Average values of the best fitness per run of scenarios I (*top graph*) and II (*bottom graph*)

results (and this is indeed the case, since Scenario III gave better results than Scenario I), and (2) the genetic operators help in mining non-local dependencies (again, note that the results for the nets with non-local non-free-choice constructs—*a6nfc* and *driversLicense*—are better in Scenario III than in Scenario II). Thus, in short, Scenario III shows that the GA was going on the right track, but it would need more iterations to reach or outperform the results of Scenario II for all nets.

- The use of genetic operators makes the population converge to better models in a faster pace, as shown in Figs. 20 and 21. Note that for scenarios III and IV, in many situations the best individuals have a fitness superior to 0.90 already at generation 100. Furthermore, the results in these figures show that the heuristics used to build the initial population indeed capture many of



**Fig. 21** Average values of the best fitness per run of scenarios III (top graph) and IV (bottom graph)

the correct causality relations. Note that the best fitness for individuals at generation 0 is already bigger than 0.0 for all nets in scenarios II and IV, and this is not the case for scenarios I and III.

- Many of the runs took on average less than 15 min (cf. Figs. 22 and 23). However, logs from nets with loops and/or parallel constructs tend to take more time. These nets usually allow for more interleaving situations and, therefore, have bigger logs. Since the fitness is calculated by replaying the logs in the individuals of a population, the genetic algorithm takes longer for these models. As a general remark, the more different traces and tasks a log contains, the higher will be the computational time of the genetic algorithm.

The next section shows the results for the experiments with noisy logs.

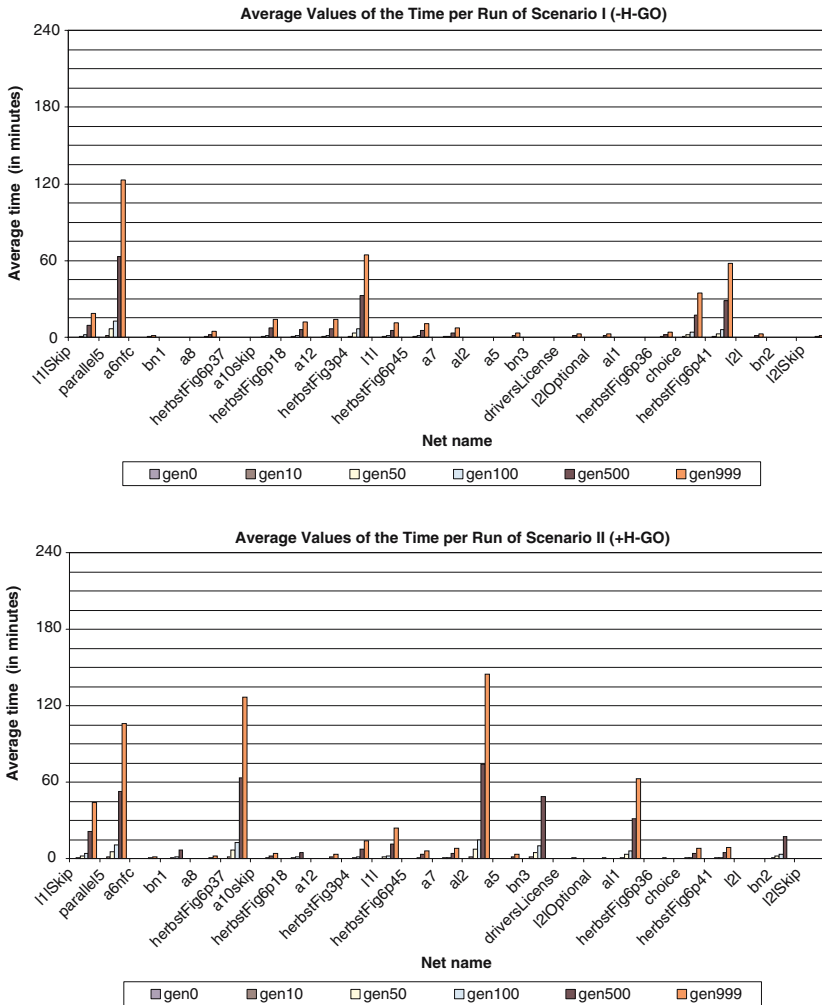
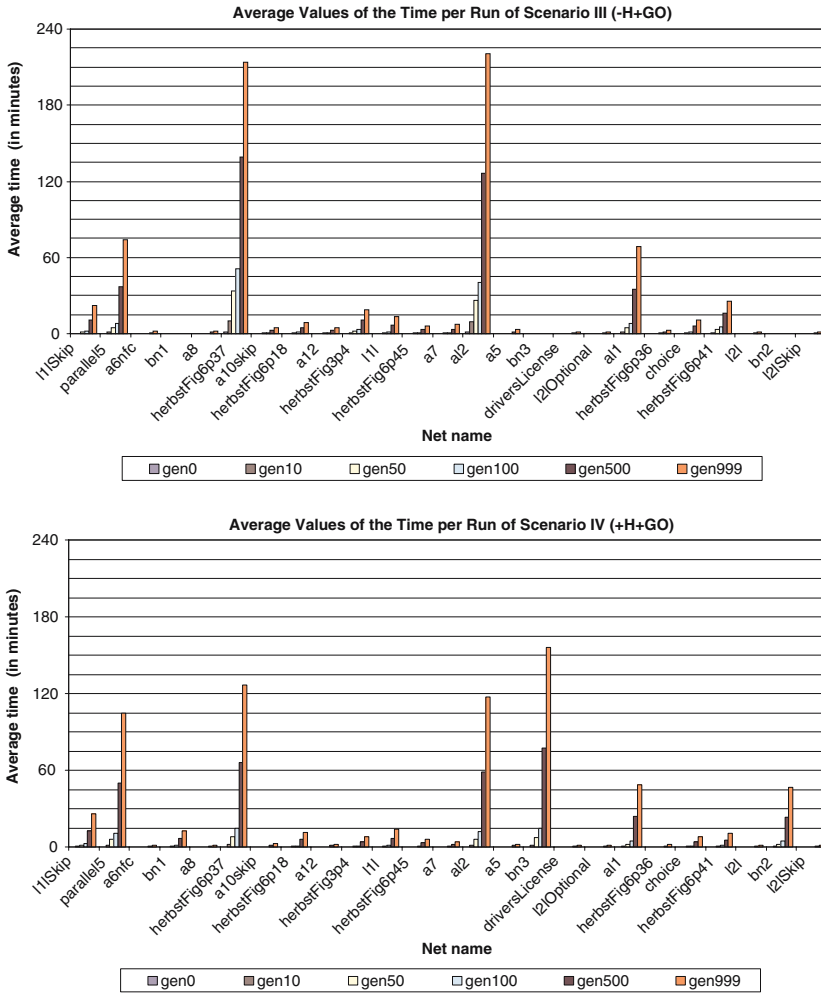


Fig. 22 Average values of the time per run of scenarios I (top graph) and II (bottom graph)

### 6.2 Experiments with noisy logs

Before reporting on the experiments with noisy logs, we explain the approach we chose to handle mined models from these logs. Noise can be defined as *low-frequent incorrect* behavior in the log. A log may contain noise because some of its traces are incomplete (e.g., they correspond to running cases in the system that have not been completed yet), or the traces reflect incorrect behavior (e.g., due to some temporal system misconfiguration). Either way, the presence of noise may hinder the correct discovery of a process model. Noisy behavior is typically difficult to detect because it cannot be easily distinguished from other *low-frequent correct* behavior in the log (for instance, the execution



**Fig. 23** Average values of the time per run of scenarios III (top graph) and IV (bottom graph)

of exceptional paths in the process). For this reason, and because the genetic algorithm is designed to always benefit individuals that can correctly parse the *most frequent* behavior in the log, we have opted for a post-processing step to “clean” mined models from the effects of noise. In short, this post-processing step works by *pruning the arcs of a (mined) model that are used fewer times than a certain threshold*.

The main advantage of a post-pruning step is that, because it works independently of the process mining algorithm, it does not avoid the discovery of low-frequent behavior. Thus, if the mined low-frequent behavior is a correct one, it can remain in the model. If the mined low-frequent behavior corresponds to noisy behavior, the end user has the possibility to clean the mined



model. Furthermore, arc post-pruning can also be used over any model to get a more concise view (in terms of the number of arcs in the model) of the most frequent behavior. As a final remark, we point out that arc post-pruning is also the approach adopted by the other related process mining techniques in Agrawal et al. (1998), Cook et al. (2004), Greco et al. (2004), and Herbst and Karagiannis (2004) to clean mined models. The remainder of this section provides more details about this arc post-pruning approach, the experiments setup and results. The aim of the experiments is to get an indication of how sensitive the genetic algorithm is to different noise types.

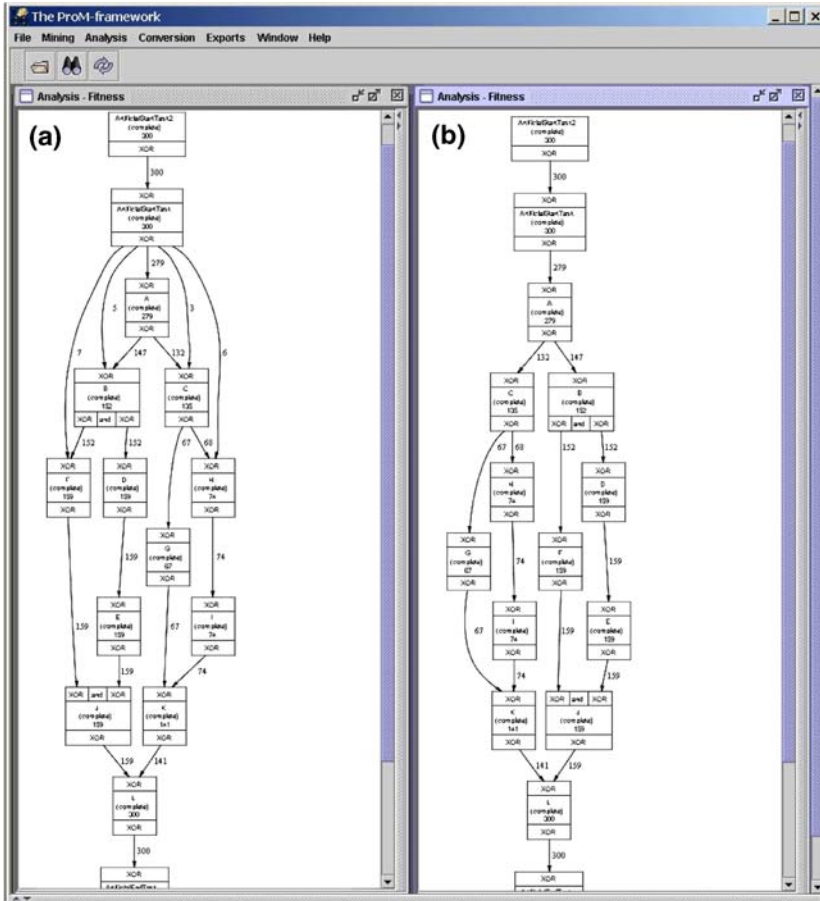
### 6.2.1 Post pruning

The post-pruning step removes the arcs that are used fewer times than a given threshold from a (mined) model. The threshold refers to the arc usage percentage. The arc usage indicates the number of times that an arc (or dependency) is used when a log is replayed by an individual. The arc usage percentage defined by the threshold is relative to the most frequently used arc. As an illustration, assume that the most frequent arc usage of a mined model to a given log is 300. If the threshold is set for 5%, all arcs of this model that are used 15 or fewer times are removed during the post-pruning step. This situation is depicted in Fig. 24. The mined model is in Fig. 24a. The pruned model is in Fig. 24b. Notice that the arcs of the mined model that were used (from left to right) 7, 5, 3, and 6 times are not shown in the pruned model. When the removal of arcs leads to dangling activities (i.e., activities without ingoing and outgoing arcs), these activities are also omitted in the post-pruned model.

### 6.2.2 Setup

The genetic algorithm was tested over noisy logs from five different process models: a12, bn1, herbstFig3p4, herbstFig6p36 and herbstFig6p37. These models<sup>18</sup> contain constructs like sequences, choices, parallelism, structured loops, and non-local non-free-choice constructs, and have between 12 and 42 tasks. The noise-free log of every net has 300 traces (actually, these are the same noise-free logs used during the experiments reported in Sect. 6.1). For every noise-free log, 12 noisy logs were generated: six logs with 5% noise and six logs with 10% noise. The six noise types used are: *missing head*, *missing body*, *missing tail*, *swap tasks*, *remove task*, and *mix all*. These noise types are the ones described in Maruster (2003). If we assume a trace  $\sigma = t_1 \dots t_{n-1} t_n$ , these noise types behave as follows. *Missing head*, *body*, and *tail*, respectively, randomly remove sub-traces of tasks in the head, body and tail of  $\sigma$ . The head goes from

<sup>18</sup> The main motivation to select only five models is that, as shown in Sect. 6.1, a run of the genetic algorithm is time consuming and, as we explain further in the text, 12 noisy logs were generated per net. Thus, the choice for five models is a pragmatic one.



**Fig. 24** Illustration of applying post-pruning to arcs of a mined model. The mined model is in (a), and the resulting post-pruned model is in (b). The numbers next to the arcs in these models inform how often these arcs have been used while replaying the log for the models. The post-pruning illustrated here has a threshold of 5%. Thus, since the highest arc usage of the mined model in (a) is 300, all of its arcs that are used 15 or fewer times are not shown in the resulting pruned model in (b)

$t_1$  to  $t_{n/3}$ .<sup>19</sup> The body goes from  $t_{(n/3)+1}$  to  $t_{(2n/3)}$ . The tail goes from  $t_{(2n/3)+1}$  to  $t_n$ . The removed sub-traces contain at least one task and at most all the tasks in the head, body or tail. *Swap task* exchanges two tasks in  $\sigma$ . *Remove task* randomly removes one task from  $\sigma$ . *Mix all* randomly performs (with the same probability) one of the other five noise types to a traces in a log. Real life logs will typically contain mixed noise. However, the separation between the noise

<sup>19</sup> The division  $n/3$  is rounded to the largest double value that is not greater than  $n/3$  and is equal to a mathematical integer.

types allows us to better assess how the different noise types affect the genetic algorithm.

To speed up the computational time of the genetic algorithm, the similar traces were grouped into a single one and a weight was added to indicate how often the trace occurs. Traces with the same sequence of tasks were grouped together. For every noisy log, 50 runs were executed. The configuration of every run is the same used for the noise-free experiments of the scenario IV (see Sect. 6.1). After a run was complete, the mined model was used as input for a post-processing step to prune its arcs. Every mined model went two post-pruning steps: one to prune with a threshold of 5% and another to prune with a threshold of 10%. So, when analyzing the results, we look at (1) the mined model returned by the genetic algorithm, (2) the model after applying 5% pruning, and (3) the model after 10% pruning. The post-processing step is implemented as the *Prune Arcs* analysis plug-in in the ProM framework.

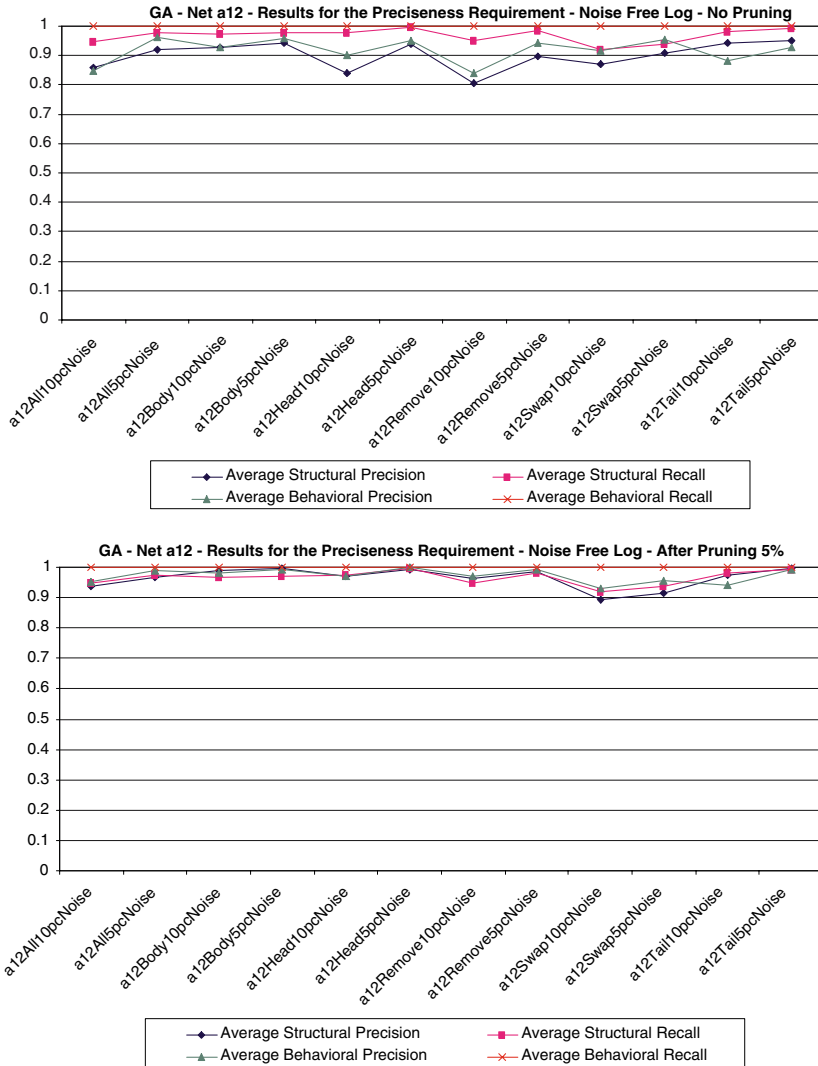
### 6.2.3 Results

The results for the experiments with logs of the net **a12** are in Figs. 25–28. We only show the results for the net **a12** because the obtained results for the other nets lead to the same conclusions that can be drawn based on the analysis of the results for **a12**. Every figure plots the results before and after pruning. However, we have omitted the results for 10% arc-pruning because the results are just like the results for 5% arc-pruning. Furthermore, for every graph, the *x*-axis shows, for a given net (or original model), the noise type and the percentage of noise in the log. For instance, **a12All10pcNoise** is a short for “Noisy log for the net **a12** (**a12**). The noise type is *mix all* (**All**) and this log contains at most 10% (**10pc**) of noise (**Noise**).” The *y*-axis contains the values for the analysis metrics.

Additionally, we have plotted the metric values for the mined model and original model with respect to the noisy log and the noise-free one. The reason is that the analysis with the noisy logs allow us to check if the mined models over-fit the data (since these noisy logs were given as input to the genetic algorithms). For instance, if some mined model can proper complete the noisy log (can parse all the traces without missing tokens or tokens left-behind), this mined model has over-fitted the data. On the other hand, when the model does not over-fit the data, the analysis with the noise-free logs can check if the mined model correctly captures the most frequent noise-free behavior (since the noisy logs used for the experiments were created by inserting noisy behavior into these noise-free logs).

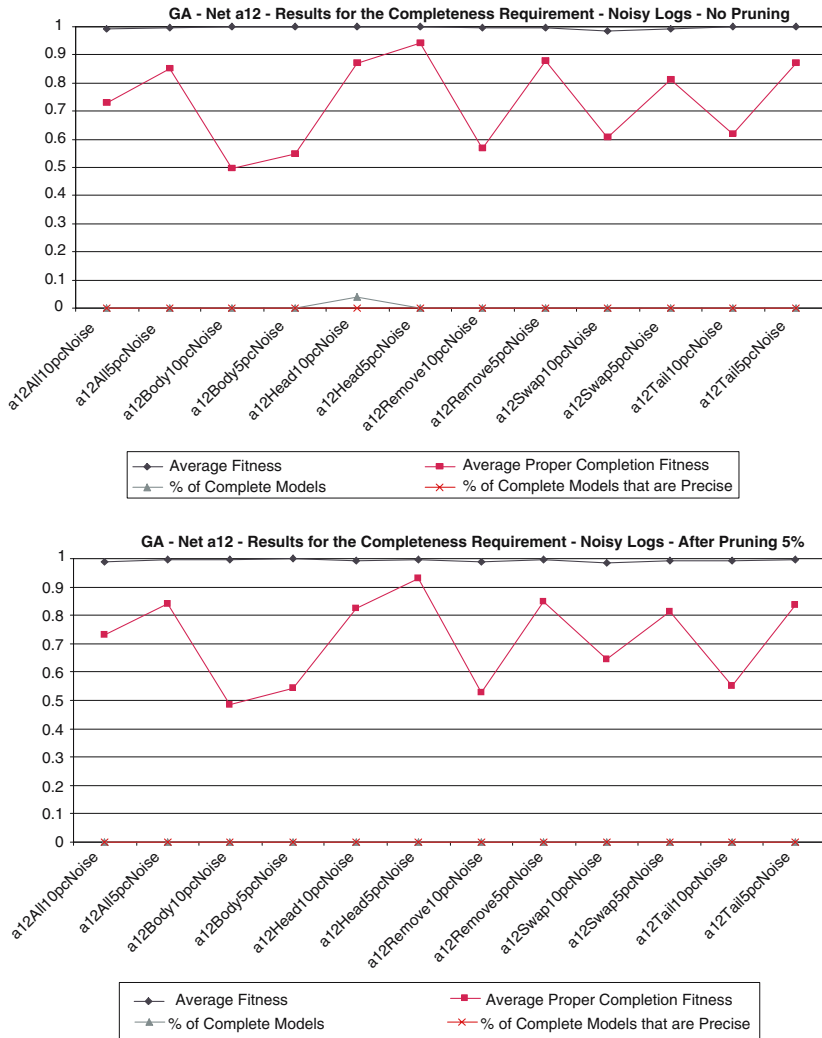
In a nutshell, we can conclude that (1) the genetic algorithm is more sensitive to the noise type *swap tasks* (and, consequently, *mix all*) and (2) the mined models do not tend to over-fit the noisy logs. More specifically, the results point out that:

- The genetic algorithm is more robust to 5% noise than to 10% noise. But the 5% arc post-pruning gave the same results as the 10% one.



**Fig. 25** Average values for the behavioral and structural precision/recall metrics of the models mined by the genetic algorithm (GA) for noisy logs of the net a12. The top (bottom) graph shows the results for the mined models before (after) arc post-pruning. The results show that (1) the mined models have a behavior that is quite similar to the original models (since behavioral precision > 0.8 and behavioral recall > 0.95), and (2) the arc post-pruning is more effective for the noise types missing head/body/tail, and remove task (since all the values plotted in the bottom graph are better than the ones in the top graph)

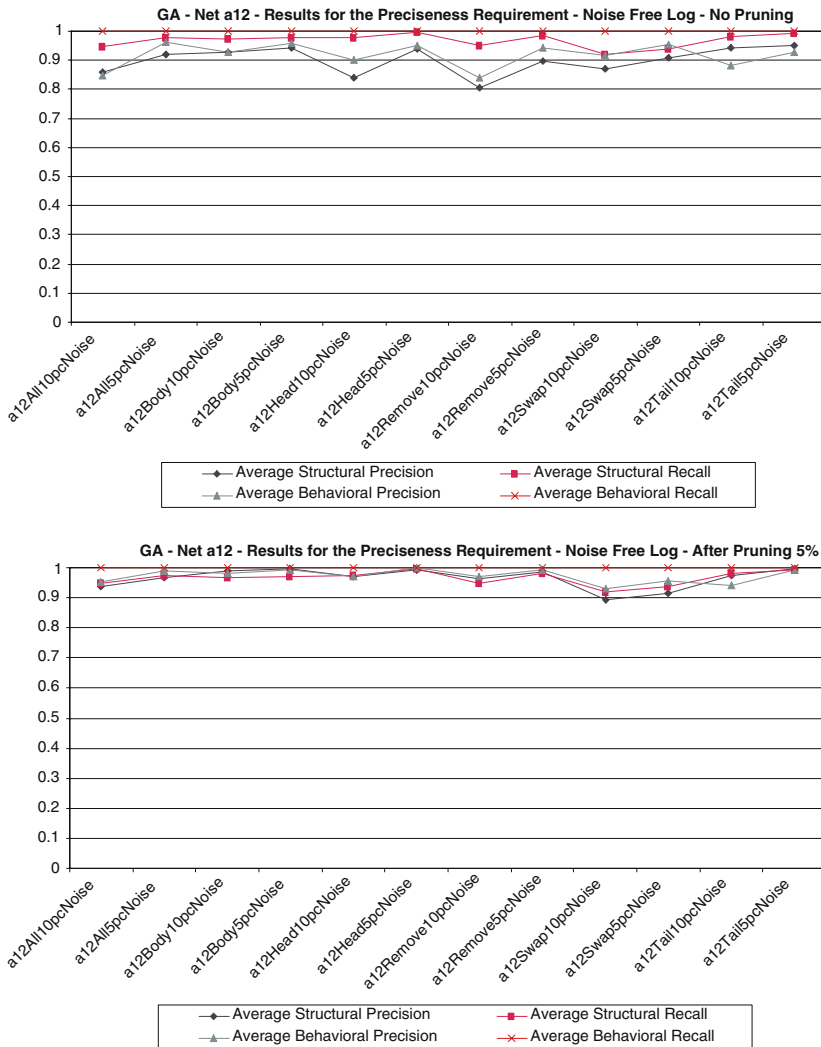
- The pruning is more effective for the noise types missing head, missing body, missing tail, and remove task. This makes sense because these noise types usually can be incorporated to the net by adding causality dependencies to skip the “removed” or “missing” tasks. In other words, the main net structure (the one also contained in the original model) does not change, only extra



**Fig. 26** Results of the GA for the noisy logs of the net a12: completeness metrics. The metrics were calculated based on the noisy logs used during the mining. The *top graph* shows the results for the mined models. The *bottom graph* shows the results after the mined models have undergone 5% arc post-pruning. Note that the top graph indicates that only two models for the noise type *missing head* (10% noise) over-fit the data. However, overall the mined models did not over-fit the data, since the average proper completion < 0.95 and the percentage of complete models is 0 for almost all logs

causality dependencies need to be added to it. This explains why the arc post-pruning works quite fine for these noise types.

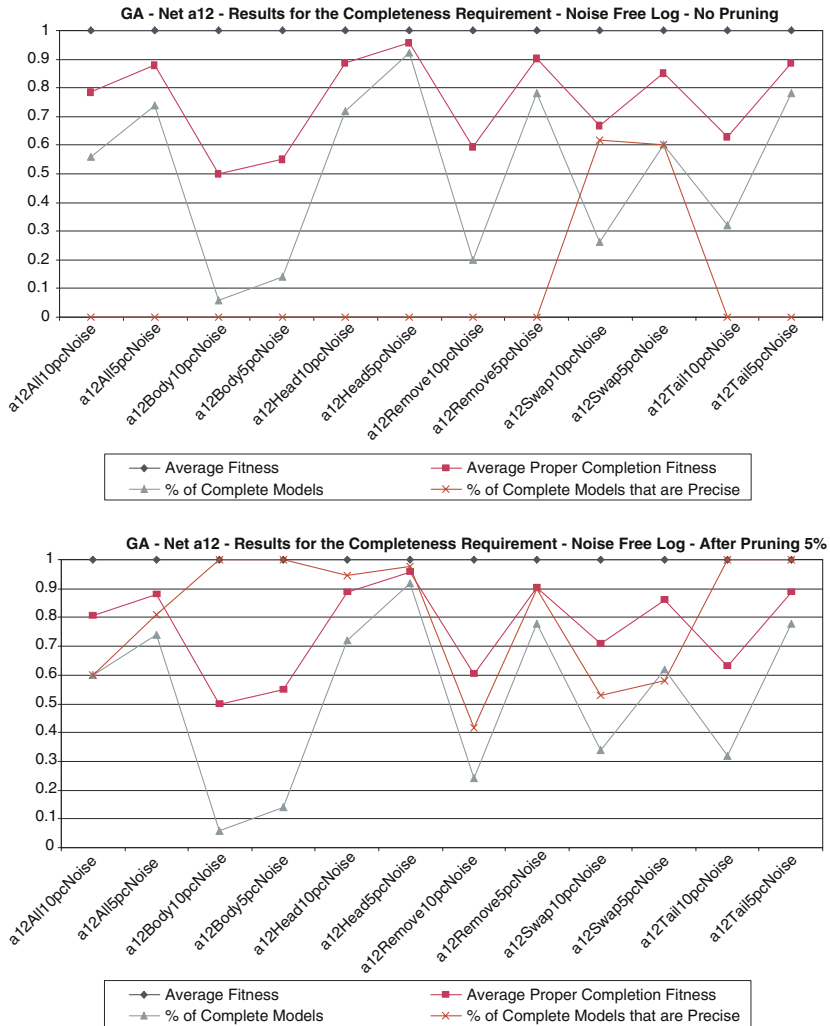
- Related to the previous item, the noise type *swap tasks* affects the quality of the mined results the most. By looking at Figs. 25 and 27, one can see that the behavioral/structural precision and recall of the mined models for logs with swapped tasks (a12Swap5pcNoise and a12Swap-10pcNoise)



**Fig. 27** The results show the same metrics as explained for Fig. 25, but these metrics are calculated base on a noise-free log of a12. Note that, contrary to the results in Fig. 25, all mined models (before and after pruning) have an average behavioral recall that is equal to 1. This means that, with respect to the noise-free log, all the behavior allowed by the original model is also captured by the mined models.

did not change dramatically after the pruning. This is probably because the over-fitting of the mined models to the logs involves more than the simple addition of causality dependencies. I.e., the main structure of mined models is more affected by the *swap tasks* noise type.

The next section shows the results of applying the hybrid version of the genetic algorithm (scenario IV) to event logs from a municipality in The Netherlands.



**Fig. 28** Same metrics as in Fig. 26, but this time the values are calculated based on a noise-free log of a12. The values for the average proper completion fitness point out that the models correctly captured at least 50% (see a12Body10pcNoise) of the behavior in the log. This indicates that the fitness indeed benefits the individuals that correctly model the most frequent behavior in the log. Besides, note that many more mined models are complete and precise after they have undergone arc post-pruning

### 7 Mining real-life logs

This section shows the results of applying the genetic algorithm for event logs from a municipality in The Netherlands. The four selected process models deal with the handling of complaints (Bezwaar, BezwaarWOZ, Afschriften) and the process to get a building permit (Bouwvergunning). As shown in Table 2, not all traces in the event logs given to us are compliant with the prescribed (original)

**Table 2** Percentage of traces that could be correctly parsed by the *original* process models

Process model	Number of traces	Correctly parsed traces (%)
Bezwaar	35	51
BezwaarWOZ	100	47
Afschriften	358	100
Bouwvergunning	407	80

A trace is correctly parsed when no tokens are missing or left behind during the parsing.

**Table 3** Percentage of traces that could be correctly parsed by the *mined* process models before and after the arc pruning

Process Model	% Correctly Parsed Traces Mined Model			
	No pruning (%)	1% pruning (%)	5% pruning (%)	10% pruning (%)
Bezwaar	100	100	62	51
BezwaarWOZ	100	95	82	82
Afschriften	100	99	99	88
Bouwvergunning	70	70	64	64

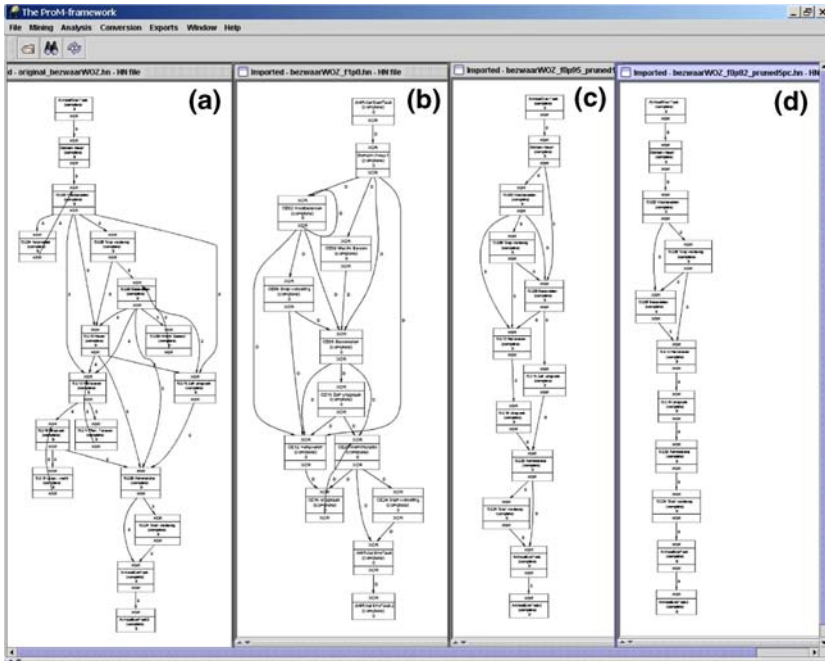
Note that process Bouwvergunning has the only log to which the mined model cannot correctly parse all the traces. This process has a construct with four tasks in parallel. As the results for the experiments in Sect. 6 indicated, the genetic algorithm goes for the most frequent behavior in these situations

process model.<sup>20</sup> So, these logs are a nice setting to test how robust the genetic algorithm is to noise.

The experiments consisted in running the algorithm for every log. The configuration used is just like the one in Sect. 6.1, but with ten individuals, at most 5,000 generations and for ten seeds only. For every log, the best mined model over all seeds was selected to undergo the arc post-pruning step (cf. Sect. 6.2). The arc pruning is a post-processing step that eliminates from a mined model the arcs that are used up to a percentage of the most frequently used arc, during the parsing of the event log by the mined model. The percentage is set by the user. The choice for a post-pruning step is based on the fact that we do not typically know how much noise an event log contains. Thus, we let the genetic algorithm mine the most frequent behavior in the log and use the post-pruning to “clean” this model even further. The results of the experiments are in Table 3. As can be seen, the algorithm indeed goes for the parsing of the most frequent behavior in the log. Note that even after pruning the arcs that are used 10 or fewer times less than the most frequently used arc, the resulting process model can still correctly parse more than half of the traces in the log. As an example, Fig. 29 shows the original and (pruned) mined models for the process BezwaarWOZ (cf. Table 3).

<sup>20</sup> All modes contain sequences, choices, length-two loops and invisible tasks. Additionally, Bouwvergunning has four tasks in parallel. The amount of tasks per process is: Afschriften (11 tasks), Bezwaar (14), BezwaarWOZ (17), and Bouwvergunning (19).





**Fig. 29** Original and mined models for the process BezwaarWOZ. The original model is in (a). The unpruned mined model is in (b). c and d, respectively, show the models after applying 1 and 5% arc pruning to the mined model in (b)

## 8 Comparison to some related approaches

The genetic algorithm is the only mining algorithm so far (cf. Sect. 9) that can mine sequence, choice, parallelism, (arbitrary) loops, invisible tasks, and non-free-choice constructs at once, while being robust to noise. Therefore, it is interesting to compare it with other existing approaches. In this section, we compare the genetic algorithm with two other mining algorithms:  $\alpha^{++}$  and HeuristicsMiner.

The aim of the comparison is to assess if these approaches outperform or not the genetic algorithm when mining logs. The reasons to select these two algorithms are: (1) both algorithms are implemented in the ProM framework (what facilitates the comparison); (2) the  $\alpha^{++}$  (Wen et al. 2006) is an extension of the  $\alpha$  algorithm in (Alves de Medeiros et al. 2004a) to also explicitly capture non-free-choice constructs in the mined model; and (3) the HeuristicsMiner<sup>21</sup> is robust to noise, can tackle invisible tasks and uses heuristics that are based on the Dependency measure explained in Sect. 4.1.1. Furthermore, both

<sup>21</sup> The HeuristicsMiner is the ProM version of the Little Thumb tool (Weijters and van der Aalst 2003).

algorithms can mine the basic structural constructs (sequence, choice, parallelism and loops).

The  $\alpha^{++}$  and HeuristicsMiner algorithms were applied<sup>22</sup> over the three sets of logs (noise-free, noisy and real-life) to which we executed the genetic algorithm. The results show that:

*Experiments with noise-free logs:* As the values for the analysis metrics (cf. Sect. 5) in Figs. 30 and 31 show, both algorithms could find a complete and precise model for many of the logs. However, as expected, the  $\alpha^{++}$  could not correctly mine models with invisible tasks (cf. Fig. 30) and the HeuristicsMiner was unable to capture non-free-choice constructs (cf. Fig. 31).

*Experiments with noisy logs:* As the  $\alpha^{++}$  is not robust to noise, none of its mined models could properly complete any of the traces in the noisy logs in Sect. 6.2. Actually, the resulting mined models were very spaghetti. However, the HeuristicsMiner outperformed the genetic algorithm for the noisy logs in Sect. 6.2. In fact, except for the model with non-free-choice (herbst-Fig6p36), the HeuristicMiner mined models with the same structure of the original models for all six noise types.

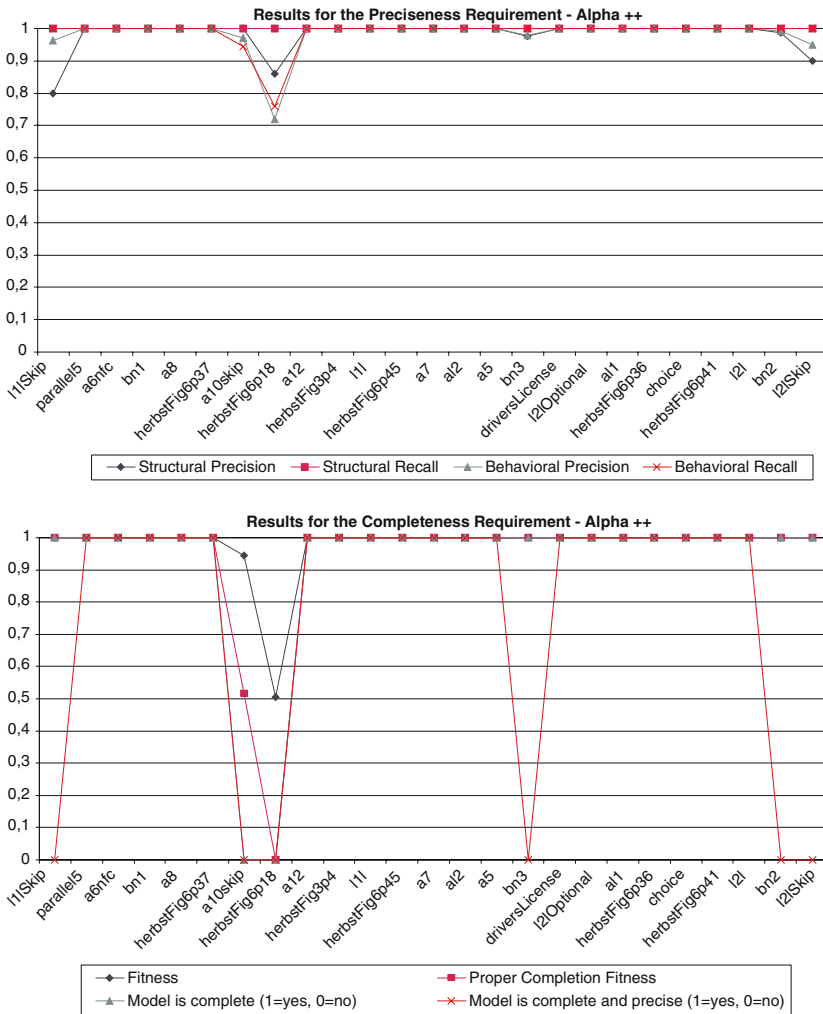
*Experiments with real-life logs:* Given the results reported in the previous two items, one would expect the HeuristicMiner to outperform the genetic algorithm when running over real-life logs (cf. Sect. 7). Especially because all the structural constructs in the prescribed models of the Dutch municipality can be tackled by the HeuristicsMiner. However, as the results in Table 4 indicate, both the  $\alpha^{++}$  and the HeuristicsMiner underperformed the genetic algorithm when correctly capturing the most frequent behavior in the log.

As a final remark, we emphasize that both the  $\alpha^{++}$  and the HeuristicsMiner are extremely fast (models were mined within 30s on the same computer used for the experiments with the genetic algorithm). However, when looking at the quality of the mined models, the genetic algorithm is better than these other two related algorithms because, for the set of logs discussed above, it mined models that correctly depict the most frequent behavior in more situations.

## 9 Related work

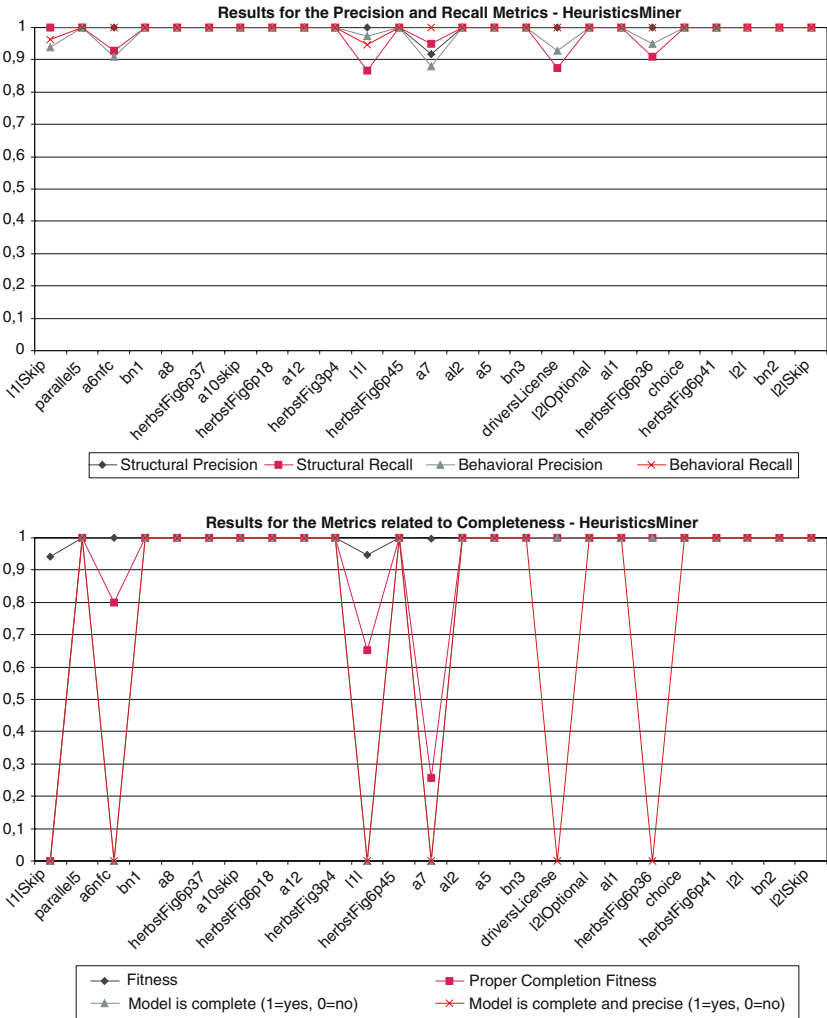
The idea of process mining is not new (van der Aalst et al. 2003; van der Aalst and Weijters 2004; Agrawal et al. 1998; Cook and Wolf 1998a, 1998b, 1999; Herbst 2000, 2001; Herbst and Karagiannis 2000, 2004; Cook et al. 2004; Maruster et al. 2002; Maxeiner et al. 2001; Schimm 2002, 2004; Pinter and Golani 2004; Weijters and van der Aalst 2003). Cook and Wolf have investigated sim-

<sup>22</sup> The default configuration provided by ProM was used when running these algorithms.



**Fig. 30** Results for the models mined by the  $\alpha^{++}$  algorithm. Note that the non-free-choice constructs are indeed correctly captured (cf. results for a6nfc, driversLicense and herbstFig6p36), but the algorithm cannot correctly mine the models with invisible tasks (cf. 111Skip, a10skip, herbstFig6p18, bn3, bn2 and l2lSkip)

ilar issues in the context of software engineering processes. In Cook and Wolf (1998a) they describe three methods for process discovery: one using neural networks, one using a purely algorithmic approach, and one Markovian approach. The authors consider the latter two the most promising approaches. The purely algorithmic approach builds a finite state machine where states are fused if their futures (in terms of possible behavior in the next  $k$ -steps) are identical. The Markovian approach uses a mixture of algorithmic and statistical methods and is able to deal with noise. Note that the results presented in Cook and Wolf



**Fig. 31** Results for the models mined by the HeuristicsMiner. Note that the invisible tasks are correctly captured, but the non-free-choice constructs are not. Furthermore, the algorithm seems to have problems with mining length-one loops (cf. 111Skip and 111)

(1998a) are limited to sequential behavior. Cook and Wolf extend their work to concurrent processes Cook and Wolf (1998b) and Cook et al. (2004). They propose specific metrics (entropy, event type counts, periodicity, and causality) and use these metrics to discover models out of event streams. However, they do not provide an approach to generate explicit process models. Recall that the final goal of the approach presented in this paper is to find explicit representations for a broad range of process models, i.e., we want to be able to generate a concrete model rather than a set of dependency relations between events.

**Table 4** Percentage of traces that could be correctly parsed by the models mined by the  $\alpha^{++}$  and the HeuristicsMiner

Process model	Correctly parsed traces mined model (%)	
	$\alpha^{++}$ (%)	HeuristicsMiner (%)
Bezwaar	0	48
BezwaarWOZ	0	86
Afschriften	100	100
Bouwver gunning	0	0

Note that both algorithm underperformed the genetic algorithm (cf. Table 3) when correctly portraying the most frequent behavior in the logs

In [Cook and Wolf \(1999\)](#) provide a measure to quantify discrepancies between a process model and the actual behavior as registered using event-based data. The idea of applying process mining in the context of workflow management was first introduced in [Agrawal et al. \(1998\)](#). This work is based on workflow graphs, which are inspired by workflow products such as IBM MQSeries workflow (formerly known as Flowmark) and InConcert. In this paper, two problems are defined. The first problem is to find a workflow graph generating events appearing in a given workflow log. The second problem is to find the definitions of edge conditions. A concrete algorithm is given for tackling the first problem. The approach is quite different from other approaches: Because the nature of workflow graphs there is no need to identify the nature (AND or XOR) of joins and splits. As shown in [Maruster et al. \(2002, Submitted\)](#) workflow graphs use true and false tokens which do not allow for cyclic graphs. Nevertheless, [Agrawal et al. \(1998\)](#) partially deals with iteration by enumerating all occurrences of a given activity and then folding the graph. However, the resulting conformal graph is not a complete model. In [Maxeiner et al. \(2001\)](#), a tool based on these algorithms is presented. [Pinter and Golani \(2004\)](#) extends the work in [Agrawal et al. \(1998\)](#) to also consider the time information in the logs and, consequently, better detect concurrent behavior. [Schimm 2002, 2004](#) has developed a mining tool suitable for discovering hierarchically structured workflow processes. This requires all splits and joins to be balanced. [Herbst \(2000, 2001\)](#) and [Herbst and Karagiannis \(2000, 2004\)](#) also address the issue of process mining in the context of workflow management using an inductive approach. The work presented in [Herbst and Karagiannis \(2000\)](#) is limited to sequential models. The approach described in [Herbst \(2001, 2000\)](#) and [Herbst and Karagiannis \(2004\)](#) also allows for concurrency. It uses stochastic activity graphs as an intermediate representation and it generates a workflow model described in the ADONIS modelling language. In the induction step activity nodes are merged and split in order to discover the underlying process. A notable difference with other approaches is that the same activity can appear multiple times in the workflow model, i.e., *the approach allows for duplicate activities*. The graph generation technique is similar to the approach of [Agrawal et al. \(1998\)](#) and [Maxeiner et al. \(2001\)](#). The nature of splits and

joins (i.e., AND or XOR) is discovered in the transformation step, where the stochastic activity graph is transformed into an ADONIS workflow model with block-structured splits and joins. In contrast to the previous papers, the work (Maruster et al. 2002; Weijters and van der Aalst 2003) is characterized by the focus on workflow processes with concurrent behavior (rather than adding ad-hoc mechanisms to capture parallelism). In Weijters and van der Aalst (2003), a heuristic approach using rather simple metrics is used to construct so-called “dependency/frequency tables” and “dependency/frequency graphs.” The preliminary results presented in Weijters and van der Aalst (2003) only provide heuristics and focus on issues such as noise. In van der Aalst and van Dongen (2002) the EMiT tool is presented which uses an extended version of the  $\alpha$ -algorithm to incorporate timing information. For a detailed description of the  $\alpha$ -algorithm and a proof of its correctness we refer to van der Aalst et al. (2005). For a detailed explanation of the constructs the  $\alpha$ -algorithm does not correctly mine and an extension to mine short-loops, (see Alves de Medeiros et al. 2003, 2004a). The main differences from our work to the above mentioned ones are that (1) our search is not primarily based on local (direct neighborhood) information in the log, (2) we can capture non-free-choice constructs, and (3) we try to discover the dependencies (causality relations) and the semantics of the split/join points all together. With respect to non-free-choice constructs, as already explained in Sect. 8, Wen et al. (2006) have extended the  $\alpha$  algorithm in Alves de Medeiros et al. (2004a) to also discover this kind of construct. The extension is called the  $\alpha^{++}$  algorithm. However, unlike the genetic algorithm, the  $\alpha^{++}$  is unable to capture invisible tasks and is also not robust to noise. The approach by Greco et al. (2006) handles non-free-choice situations by clustering the traces that follow different paths in the mined model. The difference to our approach is that the non-free-choice constructs are not explicitly captured in the control-flow structure of the mined model.

Process mining can be seen as a tool in the context of BPI. In Grigori et al. (2001) a BPI toolset on top of HP’s Process Manager is described. The BPI tools set includes a so-called “BPI Process Mining Engine.” However, this engine does not provide any techniques as discussed before. Instead it uses generic mining tools such as SAS Enterprise Miner for the generation of decision trees relating attributes of cases to information about execution paths (e.g., duration). In order to do workflow mining it is convenient to have a so-called “process data warehouse” to store audit trails. Such as data warehouse simplifies and speeds up the queries needed to derive causal relations. In Eder et al. (2002), zur Mühlen (2001) and zur Mühlen and Rosemann (2000), the design of such warehouse and related issues are discussed in the context of workflow logs. Moreover, zur Mühlen and Rosemann (2000) describes the PISA tool which can be used to extract performance metrics from workflow logs. Similar diagnostics are provided by the ARIS Process Performance Manager (PPM) (IDS Scheer 2002). The later tool is commercially available and a customized version of PPM is the Staffware Process Monitor (SPM) (Staffware 2002) which is tailored towards mining Staffware logs. Note that

none of the latter tools is extracting the process model. The main focus is on clustering and performance analysis rather than causal relations as in (Agrawal et al. 1998; Cook and Wolf 1998a,b,1999; Herbst 2000, 2001; Herbst and Karagiannis 2000; Maruster et al. 2002; Schimm 2002; Weijters and van der Aalst 2003).

More from a theoretical point of view, the rediscovery problem discussed in this paper is related to the work discussed in Angluin and Smith (1983), Gold (1978) and Pitt (1889). In these papers the limits of inductive inference are explored. For example, in Gold (1978) it is shown that the computational problem of finding a minimum finite-state acceptor compatible with given data is NP-hard. Several of the more generic concepts discussed in these papers could be translated to the domain of process mining. It is possible to interpret the problem described in this paper as an inductive inference problem specified in terms of rules, a hypothesis space, examples, and criteria for successful inference. The comparison with literature in this domain raises interesting questions for process mining, e.g., how to deal with negative examples (i.e., suppose that besides  $\log W$  there is a  $\log V$  of traces that are not possible, e.g., added by a domain expert). However, despite the many relations with the work described in Angluin and Smith (1983), Gold (1978) and Pitt (1889) there are also many differences, e.g., we are mining at the net level rather than sequential or lower level representations (e.g., Markov chains, finite state machines, or regular expressions). For a survey of existing research, we also refer to van der Aalst et al. (2003).

There have been some papers combining Petri nets and genetic algorithms, (cf. Bourdeaud'huy and Yim 2002; Malpathak et al. 2002; Mauch 2003; Moore and Hahn 2003a,b, 2004; Nummela and Julstrom 2005; Reddy et al. 2001; Tohme et al. 1999). However, these papers do not try to discover a process model based on some event log. The approach in this paper is the first approach using genetic algorithms for process discovery. The goal of using genetic algorithms is to tackle problems such as duplicate activities, hidden activities, non-free-choice constructs, noise, and incompleteness, i.e., overcome the problems of some of the traditional approaches. Actually, we have also previous papers on genetic process mining (Alves de Medeiros et al. 2006, van der Aalst et al. 2005). However, the work presented in this paper differs from our previous papers for the following reasons. In van der Aalst et al. (2005), the causal matrix representation (cf. Definition 5) was more restrictive because the subsets in the input and output condition functions (i.e.,  $I$  and  $O$ ) were *partition sets* of  $A$ . Thus, nets like the one in Figure 5 could not be supported. The work in Alves de Medeiros et al. (2006) removed these restrictions and improved the fitness measure. However, neither the fitness measure in van der Aalst et al. (2005) nor the one in Alves de Medeiros et al. (2006) have the preciseness requirement (cf. Sect. 4.2.2) to punish over-general models. Additionally, the works in Alves de Medeiros et al. (2006) and van der Aalst et al. (2005) (1) do not include the analysis metrics (cf. Sect. 5) to assess the quality of the mined models, (2) have experiments with a smaller set of synthetic logs, and (3) do not include experiments with real-life logs.

With respect to the analysis metrics, some notions related to the behavioral precision metric (cf. Definition 13) have been used in Greco et al. (2006) and Rozinat and van der Aalst (2005). In Greco et al. (2006), the notion of *soundness* is used to check for how much extra behavior a mined model allows for. However, the metric assumes that the target model does not contain loops. In Rozinat and van der Aalst (2005), the notion of *behavioral appropriateness* is defined to also check if a model can generate more behavior than the one expressed in the log. The problem here is that the metric cannot state when a model is precise enough. Because the analysis metrics defined in Sect. 5 make use of the original model as well, they can precisely quantify how complete and precise a mined model is.

## 10 Conclusions and future work

In this paper, we have presented a hybrid genetic algorithm to mine process models from event logs. The internal representation (the causal matrix) supports more complex routing constructs like non-free choice and invisible task. The fitness measure benefits the individuals that are complete (can parse most of the behavior in the log) and precise (cannot parse more behavior than the one that can be derived from the log). The genetic operators (crossover and mutation) manipulate the basic genetic material in the algorithm: the causality relations. The experiments with synthetic logs show that (1) the use of heuristics is indeed beneficial to speed the search performed by the genetic algorithm and (2) the genetic operators are playing their role in finding the *non-local* causality relations that can never be introduced by the heuristics. Furthermore, the experiments with synthetic noisy logs and real-life logs show that the genetic algorithm is capturing the most frequent behavior in the log even in the presence of noise. However, although the GA is able to mine models with all structural constructs but duplicates tasks and is robust to noise, it has a drawback that cannot be neglected: the computational time. For this reason, future work will focus on developing better strategies to perform this search. Additionally, we want to develop a genetic algorithm that can also mine process models with duplicate tasks. As a final remark, the genetic algorithm, analysis metrics and arc post-pruning step are all implemented as plug-ins in the ProM framework ([www.processmining.org](http://www.processmining.org)).



**Appendix A: details about the models for experiments with synthetic logs**

**Table 5** Overview of the structural constructs contained in the models for the experiments with synthetic logs

Net	Sequence	Choice	Parallelism	Length-one loop	Length-two loop	Structured loop	Arbitrary loop	Non-local NFC	Invisible tasks	Duplicates in sequence	Duplicates in parallel	Number of visible tasks
a10skip	✓	✓	✓						✓			12
a12	✓	✓	✓									14
a5	✓	✓	✓	✓								7
a6nfc	✓	✓	✓					✓				8
a7	✓	✓	✓									8
a8	✓	✓	✓									10
a1	✓	✓	✓				✓					9
a2	✓	✓	✓				✓					13
bn1	✓	✓	✓						✓			42
bn2	✓	✓	✓			✓			✓			42
bn3	✓	✓	✓			✓			✓			42
choice	✓	✓	✓									12
driversLicense	✓	✓	✓									9
herbstFig3p4	✓	✓	✓			✓		✓				12
herbstFig6p18	✓	✓	✓			✓			✓			7
herbstFig6p36	✓	✓	✓	✓	✓							12
herbstFig6p37	✓	✓	✓									16
herbstFig6p41	✓	✓	✓									16
herbstFig6p45	✓	✓	✓									8
11	✓	✓	✓	✓								6
11Skip	✓	✓	✓	✓					✓			6
12	✓	✓	✓		✓							6
12Optional	✓	✓	✓		✓							6
12ISkip	✓	✓	✓		✓							6
parallel5	✓	✓	✓		✓				✓			10

## References

- van der Aalst WMP, Alves de Medeiros AK, Weijters AJMM (2005) Genetic process mining. In: Proceedings of the 26th international conference on applications and theory of Petri nets. Lecture notes in computer science, vol 3536. Springer, Miami
- van der Aalst WMP, van Dongen BF (2002) Discovering workflow performance models from timed logs. In: Han Y, Tai S, Wikarski D (eds) International conference on engineering and deployment of cooperative information systems (EDCIS 2002). Lecture notes in computer science, vol 2480. Springer, Berlin, pp 45–63
- van der Aalst WMP, van Dongen BF, Herbst J, Maruster L, Schimm G, Weijters AJMM (2003) Workflow mining: a survey of issues and approaches. *Data Knowl Eng* 47(2):237–267
- van der Aalst WMP, Song M (2004) Mining social networks: uncovering interaction patterns in business processes. In: Desel J, Pernici B, Weske M (eds) International conference on business process management (BPM 2004). Lecture notes in computer science, vol 3080. Springer, Berlin, pp 244–260
- van der Aalst WMP, Weijters AJMM (eds) (2004) Process mining. Special issue of computers in industry, vol 53. Elsevier, Amsterdam
- van der Aalst WMP, Weijters AJMM, Maruster L (2004) Workflow mining: discovering process models from event logs. *IEEE Trans Knowl Data Eng* 16(9):1128–1142
- Agrawal R, Gunopulos D, Leymann F (1998) Mining process models from workflow logs. In: Ramos I, Alonso G, Schek H-J, Saltor F (eds) Advances in database technology—EDBT'98: 6th international conference on extending database technology. Lecture notes in computer science, vol 1377. Springer-Verlag, London, UK, pp 469–483 (ISBN: 3-540-64264-1)
- Alves de Medeiros AK, van Dongen BF, van der Aalst WMP, Weijters AJMM (2004a) Process mining: extending the  $\alpha$ -algorithm to mine short loops. BETA working paper series, WP 113, Eindhoven University of Technology, Eindhoven
- Alves de Medeiros AK, Weijters AJMM, van der Aalst WMP (2004b) Using genetic algorithms to mine process models: representation, operators and results. BETA working paper series, WP 124, Eindhoven University of Technology, Eindhoven
- Alves de Medeiros AK, Weijters AJMM, van der Aalst WMP (2006) Genetic process mining: a basic approach and its challenges. In: Business process management 2005 workshops. Lecture notes in computer science, vol 3812. Springer, Berlin, pp 203–215
- Alves de Medeiros AK, van der Aalst WMP, Weijters AJMM (2003) Workflow mining: current status and future directions. In: Meersman R, Tari Z, Schmidt DC (eds) On the move to meaningful internet systems 2003: CoopIS, DOA, and ODBASE. Lecture notes in computer science, vol 2888. Springer, Berlin, pp 389–406
- Angluin D, Smith CH (1983) Inductive inference: theory and methods. *Comput Surv* 15(3):237–269
- Bourdeaud'huy T, Yim P (2002) Petri net controller synthesis using genetic search. In: Proceedings of the 2nd IEEE international conference on systems, man and cybernetics (SMC'02), vol 1, IEEE Computer Society Press, Hammamet, Tunisia, 6–9 October 2002, pp 528–533
- Cook JE, Du Z, Liu C, Wolf AL (2004) Discovering models of behavior for concurrent workflows. *Comput Ind* 53(3):297–319
- Cook JE, Wolf AL (1998b) Event-based detection of concurrency. In: Proceedings of the 6th international symposium on the foundations of software engineering (FSE-6). ACM Press, New York, NY, USA, pp 35–45
- Cook JE, Wolf AL (1999) Software process validation: quantitatively measuring the correspondence of a process to a model. *ACM Trans Softw Eng Methodol* 8(2):147–176
- Cook JE, Wolf AL (1998a) Discovering models of software processes from event-based data. *ACM Trans Softw Eng Methodol* 7(3):215–249
- Desel J, Esparza J (1995) Free choice Petri nets. Cambridge tracts in theoretical computer science, vol 40. Cambridge University Press, Cambridge, UK
- Dumas M, van der Aalst WMP, ter Hofstede AH (eds) (2005) Process-aware information systems: bridging people and software through process technology. Wiley, New York
- Eder J, Olivotto GE, Gruber W (2002) A data warehouse for workflow logs. In: Han Y, Tai S, Wikarski D (eds) International conference on engineering and deployment of cooperative information systems (EDCIS 2002). Lecture notes in computer science, vol 2480. Springer, Berlin, pp 1–15

- Eiben AE, Smith JE (2003) Introduction to evolutionary computing. Natural computing. Springer, Berlin
- van Glabbeek RJ, Weijland WP (1996) Branching time and abstraction in bisimulation semantics. *J ACM* 43(3):555–600
- Gold EM (1978) Complexity of automaton identification from given data. *Inform Control* 37(3):302–320
- Greco G, Guzzo A, Pontieri L (2005) Mining hierarchies of models: from abstract views to concrete specifications. In: van der Aalst WMP, Benatallah B, Casati F, Curbera F (eds) Business process management. Lectures notes in computer science, vol 3649. Springer-Verlag, Berlin, Nancy, France, 5–8 September, 2005, pp 32–47
- Greco G, Guzzo A, Pontieri L, Sacca D (2004) Mining expressive process models by clustering workflow traces. In: Dai H, Srikant R, Zhang C (eds) PAKDD. Lecture notes in computer science, vol 3056. Springer, Berlin, pp 52–62
- Greco G, Guzzo A, Pontieri L, Sacca D (2006) Discovering expressive process models by clustering log traces. *IEEE Trans Knowl Data Eng* 18(8):1010–1027
- Grigori D, Casati F, Dayal U, Shan MC (2001) Improving business process quality through exception understanding, prediction, and prevention. In: Apers P, Atzeni P, Ceri S, Paraboschi S, Ramamohanarao K, Snodgrass R (eds) Proceedings of 27th international conference on very large data Bases (VLDB'01). Morgan Kaufmann, Los Atlos, CA, pp 159–168
- Grunwald PD, Myung IJ, Pitt M (eds) (2005) Advances in minimum description length theory and applications. MIT Press, Cambridge, MA
- Herbst J (2000) Dealing with concurrency in workflow induction. In: Baake U, Zobel R, Al-Akaidi M (eds) European concurrent engineering conference. SCS, Europe
- Herbst J (2001) Ein induktiver Ansatz zur Akquisition und Adaption von Workflow-Modellen. Ph.D. thesis, Universität Ulm
- Herbst J, Karagiannis D (2000) Integrating machine learning and workflow management to support acquisition and adaptation of workflow models. *Int J Intell Syst Account Finance Manag* 9:67–92
- Herbst J, Karagiannis D (2004) Workflow mining with InWoLvE. *Comput Ind* 53(3):245–264
- IDS Scheer (2002) ARIS process performance manager (ARIS PPM). <http://www.ids-scheer.com>
- Malpathak S, Saitou K, Qvam H (2002) Robust design of flexible manufacturing systems using colored Petri net and genetic algorithm. *J Int Manufact* 13(5):339–351
- Maruster L (2003) A machine learning approach to understand business processes. Ph.D. thesis, Eindhoven University of Technology, Eindhoven, The Netherlands
- Maruster L, Weijters AJMM, van der Aalst WMP, van den Bosch A (2002) Process mining: discovering direct successors in process logs. In: Proceedings of the 5th international conference on discovery science (discovery science 2002). Lecture notes in artificial intelligence, vol 2534. Springer, Berlin, pp 364–373
- Mauch H (2003) Evolving Petri nets with a genetic algorithm. In: Cantú-Paz E, Foster JA, Deb K, Davis L, Roy R, O'Reilly U, Beyer H, Standish RK, Kendall G, Wilson SW, Harman M, Wegener J, Dasgupta D, Potter MA, Schultz AC, Dowsland KA, Jonoska N, Miller JF (eds) Genetic and evolutionary computation—GECCO 2003, genetic and evolutionary computation conference, Chicago, IL, USA, 12–16 July 2003. Proceedings, Part II. Lecture notes in computer science, vol 2724. Springer, Berlin, pp 1810–1811
- Maxeiner MK, Küspert K, Leymann F (2001) Data mining von workflow-protokollen zur teilautomatisierten konstruktion von prozessmodellen. In: Proceedings of datenbanksysteme in Büro, technik und Wissenschaft. Informatik Aktuell Springer, Berlin, Germany, pp 75–84
- Milner R, Parrow J, Walker D (1992) A calculus of mobile processes. *Inform Comput* 100(1):1–77
- Moore JH, Hahn LW (2004) An improved grammatical evolution strategy for hierarchical Petri net modeling of complex genetic systems. In: Raidl GR et al (eds) Applications of evolutionary computing, Evo Workshops 2004. Lecture notes in computer science, vol 3005. Springer, Berlin, pp 63–72
- Moore JH, Hahn LW (2003a) Grammatical evolution for the discovery of Petri net models of complex genetic systems. In: Cantú-Paz E, Foster JA, Deb K, Davis L, Roy R, O'Reilly U, Beyer H, Standish RK, Kendall G, Wilson SW, Harman M, Wegener J, Dasgupta D, Potter MA, Schultz AC, Dowsland KA, Jonoska N, Miller JF (eds) Genetic and evolutionary computation—GECCO 2003, genetic and evolutionary computation conference, Chicago, IL, USA, 12–16 July 2003. Proceedings, Part II. Lecture notes in computer science, vol 2724. Springer, Berlin, pp 2412–2413.

- Moore JH, Hahn LW (2003b) Petri net modeling of high-order genetic systems using grammatical evolution. *BioSystems* 72(2):177–186.
- zur Mühlen M (2001) Process-driven management information systems combining data warehouses and workflow technology. In: Gavish B (ed) *Proceedings of the international conference on electronic commerce research (ICECR-4)*. IEEE Computer Society Press, Los Alamitos, CA, pp 550–566
- zur Mühlen M, Rosemann M (2000) Workflow-based process monitoring and controlling—technical and organizational issues. In: Sprague R (ed) *Proceedings of the 33rd Hawaii international conference on system science (HICSS-33)*. IEEE Computer Society Press, Los Alamitos, CA, pp 1–10
- Murata T (1989) Petri nets: properties, analysis and applications. *Proc IEEE* 77(4):541–580
- Nummela J, Julstrom BA (2005) Evolving Petri nets to represent metabolic pathways. In: Beyer H, O'Reilly U (eds) *GECCO*. ACM, New York, pp 2133–2139
- Pinter SS, Golani M (2004) Discovering workflow models from activities lifespans. *Comput Ind* 53(3):283–296
- Pitt L (1889) Inductive inference, DFAs, and computational complexity. In: Jantke KP (ed) *Proceedings of international workshop on analogical and inductive inference (AII)*. Lecture notes in computer science, vol 397. Springer, Berlin, pp 18–44
- Reddy JP, Kumanan S, Chetty OVK (2001) Application of Petri nets and a genetic algorithm to multi-mode multi-resource constrained project scheduling. *Int J Adv Manufact Technol* 17(4):305–314
- Reisig W, Rozenberg G (ed) (1998) *Lectures on Petri nets I: basic models*. Lecture notes in computer science, vol 1491. Springer, Berlin
- Rozinat A, van der Aalst WMP (2005) Conformance testing: measuring the fit and appropriateness of event logs and process models. In: Bussler C, Haller A (eds) *Business process management workshops*. Lecture notes in computer science, vol 3812. Springer-Verlag, Berlin, pp 163–176
- Schimm G. Process mining. <http://www.processmining.de/>
- Schimm G (2002) Process miner—a tool for mining process schemes from event-based data. In: Flesca S, Ianni G (eds) *Proceedings of the 8th European conference on artificial intelligence (JELIA)*. Lecture notes in computer science, vol 2424. Springer, Berlin, pp 525–528
- Schimm G (2004) Mining exact models of concurrent workflows. *Comput Ind* 53(3):265–281
- Staffware (2002) Staffware process monitor (SPM). <http://www.staffware.com>
- Tohme H, Nakamura M, Hachiman E, Onaga K (1999) Evolutionary Petri net approach to periodic job-shop-scheduling. In: *Proceedings of the IEEE international conference on systems, man, and cybernetics*, vol 4, pp 441–446
- Weijters AJMM, van der Aalst WMP (2003) Rediscovering workflow models from event-based data using little thumb. *Integr Comput Aided Eng* 10(2):151–162
- Wen L, Wang J, Sun J (2006) Detecting implicit dependencies between tasks from event logs. In: Zhou X, Li J, Shen HT, Kitsuregawa M, Zhang Y (eds) *APWeb*. Lecture notes in computer science, vol 3841. Springer, Berlin, pp 591–603