

## Correlation miner

***Citation for published version (APA):***

Pourmirza, S., Dijkman, R. M., & Grefen, P. W. P. J. (2017). Correlation miner: mining business process models and causal relations without case identifiers. *International Journal of Cooperative Information Systems*, 26(2), Article 1742002. <https://doi.org/10.1142/S0218843017420023>

***DOI:***

[10.1142/S0218843017420023](https://doi.org/10.1142/S0218843017420023)

***Document status and date:***

Published: 01/06/2017

***Document Version:***

Accepted manuscript including changes made at the peer-review stage

***Please check the document version of this publication:***

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

***General rights***

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.tue.nl/taverne](http://www.tue.nl/taverne)

***Take down policy***

If you believe that this document breaches copyright please contact us at:

[openaccess@tue.nl](mailto:openaccess@tue.nl)

providing details and we will investigate your claim.

International Journal of Cooperative Information Systems  
© World Scientific Publishing Company

## CORRELATION MINER: MINING BUSINESS PROCESS MODELS AND CAUSAL RELATIONS WITHOUT CASE IDENTIFIERS

SHAYA POURMIRZA, REMCO DIJKMAN, AND PAUL GREFEN  
*School of Industrial Engineering, Eindhoven University of Technology*  
*Den Dolech 2, P.O. Box 513, 5600MB*  
*Eindhoven, The Netherlands*  
*{s.pourmirza, r.m.dijkman, p.w.p.j.grefen}@tue.nl*

Received (Day Month Year)

Revised (Day Month Year)

Process discovery algorithms aim to capture process models from event logs. These algorithms have been designed for logs in which events that belong to the same case are related to each other – and to that case – by means of a unique case identifier. However, in service-oriented systems these case identifiers are rarely stored beyond request-response pairs, which makes it hard to relate events that belong to the same case. This is known as the *correlation challenge*. This paper addresses the correlation challenge by introducing a technique, called the *correlation miner*, that facilitates discovery of business process models when events are not associated with a case identifier. It extends previous work on the correlation miner, by not only enabling the discovery of the process model, but also detecting which events belong to the same case. Experiments performed on both synthetic and real-world event logs show the applicability of the correlation miner. The resulting technique enables us to observe a service-oriented system and determine – with high accuracy – which request-response pairs sent by different communicating parties are causally related to each other.

*Keywords:* Process Mining, Process Discovery, Event Correlation, Business Process Intelligence (BPI), Business Process Management (BPM)

### 1. Introduction

Over the past decade, there has been an increasing interest in the area of process mining<sup>1,2,3,4,5,6,7,8</sup>. The goal of process mining is to extract information about processes from event logs, i.e., execution histories. One of the prominent branches of process mining is process discovery<sup>1</sup>, which concerns itself with generating a business process model from an event log.

Process discovery techniques assume that an event log contains at least, for each recorded event: (i) a reference to the executed activity, (ii) a reference to the case for which the activity was executed, and (iii) the timestamp at which the activity was completed<sup>2</sup>. Table 1 shows an example of the event log involving 30 events in 10 cases. The main idea behind discovery algorithms is to merge, cluster and aggregate the different cases in an event log and generate a suitable business process model

2 *S. Pourmirza, R. Dijkman & P. Grefen*

based on that. For example, considering the event log from Table 1, it is possible to capture three different traces of execution:  $\langle A, B, E \rangle$  for Cases 1, 4 and 5;  $\langle A, D, E \rangle$  for Cases 2, 3, 8 and 9; and finally,  $\langle A, C, E \rangle$  for Cases 6, 7 and 10. Fig. 1 presents the corresponding business process model, as it would have been mined by the Disco process mining tool <sup>9</sup>.

Table 1: An example event log

Case	Activity	Timestamp	Case	Activity	Timestamp	Case	Activity	Timestamp
1	A	00:20	4	B	05:04	7	E	09:17
1	B	02:04	4	E	07:26	8	A	06:20
1	E	02:32	5	A	03:40	8	D	08:36
2	A	02:15	5	B	05:59	8	E	10:03
2	D	03:14	5	E	07:49	9	A	06:41
2	E	05:06	6	A	04:18	9	D	08:56
3	A	02:27	6	C	07:08	9	E	10:20
3	D	04:17	6	E	09:05	10	A	07:13
3	E	06:51	7	A	05:54	10	C	09:10
4	A	03:06	7	C	07:30	10	E	10:26

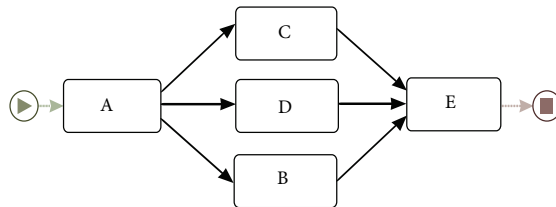


Fig. 1: Business process model generated from the example event log

While process mining techniques require an event to be associated with a case, making this association can be problematic. This is also referred to as the *correlation challenge* <sup>5</sup>. This challenge arises, for example, when information about the occurrence of activities is stored in separate service-oriented systems and there is no obvious identifier to correlate the occurrences. It is considered especially problematic in service mining <sup>10</sup>, because correlations between messages that are being exchanged for the same case, may not be stored beyond request-response pairs. When associations between events that belong to the same case are not present, process discovery becomes impossible, as illustrated by Table 2, which shows the event log from Table 1 without case identifiers. In this event log it is not possible

to identify cases and, consequently, it is not possible to cluster and aggregate them into a business process model. Moreover, we cannot see which events belong to the same case. For example, both event *A* that happens at 02:15 and event *A* that happens at 02:27 could be related to *D* that happens at 03:14.

Table 2: An example event log without case identifiers

Activity	Timestamp	Activity	Timestamp	Activity	Timestamp
<i>A</i>	00:20	<i>B</i>	05:04	<i>C</i>	07:30
<i>B</i>	02:04	<i>E</i>	05:06	<i>E</i>	07:49
<i>A</i>	02:15	<i>A</i>	05:54	<i>D</i>	08:36
<i>A</i>	02:27	<i>B</i>	05:59	<i>D</i>	08:56
<i>E</i>	02:32	<i>A</i>	06:20	<i>E</i>	09:05
<i>A</i>	03:06	<i>A</i>	06:41	<i>C</i>	09:10
<i>D</i>	03:14	<i>E</i>	06:51	<i>E</i>	09:17
<i>A</i>	03:40	<i>C</i>	07:08	<i>E</i>	10:03
<i>D</i>	04:17	<i>A</i>	07:13	<i>E</i>	10:20
<i>A</i>	04:18	<i>E</i>	07:26	<i>E</i>	10:26

Therefore, the goal and contribution of this paper is to develop a process discovery technique, correlation miner, that can:

Phase I discover business process models when events are neither correlated by case identifiers nor by additional data elements; and

Phase II discover the causal relations between the events themselves, which determine the events that belong to the same case.

Note that, this paper is an extension of our previous paper <sup>11</sup>, which presents only the first phase of the correlation miner. It extends this work by also enabling the discovery of causal relations between events in an event log. This extension is mainly presented in Section 5 and 6.

Although extensive research has been carried out on both process discovery and event correlation techniques, to the best of our knowledge, only one <sup>12</sup> such process discovery technique exists that will be further discussed in this paper. Additionally, event correlation techniques do exist, but they require additional data elements (such as customer names) to perform the correlation on.

The remainder of this paper is structured as follows. Section 2 presents some preliminaries regarding the event logs and business process models. Section 3 introduces the first phase of correlation miner that is used to capture business process models from event logs without case identifiers and Section 4 discusses the evaluation of this phase. Section 5 introduces the second phase of correlation miner that is employed to reconstruct cases from event logs without case identifiers, and subsequently, Section 6 presents the evaluation of this phase. Section 7 compares

4 *S. Pourmirza, R. Dijkman & P. Grefen*

our study with related literatures. Finally, Section 8 concludes the paper by giving a brief overview and its findings.

## 2. Preliminaries

In the literature, an event log has been defined as a multiset of cases <sup>5</sup>. However in this paper, we cannot reuse this definition since we have no explicit case. Therefore, we define an event log as follows.

Let  $A$  be a set of activities and  $T$  be a set of timestamps.  $L \subseteq A \times T$  is an *event log*. We use  $E_a$  to denote a set of events referring to the activity  $a \in A$ , such that  $E_a = \{(a, t) | (a, t) \in L\}$ . In addition, we represent the timestamp at which event  $e \in E_a$  has happened by  $t_e$ . Note that, strictly speaking, multiple events for the same activity can happen at the same time. However, for simplicity, this is not covered by our definition. While the occurrence of identical events is sufficiently rare to not create problems with the algorithm, it does require some pre-processing of the log to make sure that identical events are distinguished from each other (e.g. by adding a suffix to the timestamp).

There exist various notations in which the business process model that is the result of applying a process discovery algorithm can be represented (e.g. Petri-Nets, BPMN and YAWL) <sup>5</sup>. In this paper, we adopt a rather abstract definition of a business process model, which is in line with the notation that is used by the Disco tool for process discovery <sup>9</sup>.

The business process model is a directed *weighted* graph  $G = (\mathcal{V}, \mathcal{E}, \omega)$  such that:

- $\mathcal{V} = \{(a, n) | a \in A \wedge n = |E_a|\}$ , where  $A$  is the set of activities and  $n$  indicates for each node the number of occurrences of that activity in the event log;
- $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  is the set of edges;
- $\omega : \mathcal{E} \rightarrow \mathbb{N}$  maps edges to natural numbers, where  $\omega(a, b) = n$  if and only if there are  $n$  observations of activity  $a$  being directly followed by activity  $b$  in some case in the log. We require that  $e \in \mathcal{E}$  if and only if  $\omega(e) > 0$ .

For a node  $v \in \mathcal{V}$ , the in-degree, denoted  $deg^{in}(v)$ , is  $\sum_{(w,v) \in \mathcal{E}} \omega(w, v)$ . The out-degree, denoted  $deg^{out}(v)$ , is  $\sum_{(v,w) \in \mathcal{E}} \omega(v, w)$ . When  $deg^{in}(v) = 0$ , we call the node a source node. When  $deg^{out}(v) = 0$ , we call the node a sink node.

It is important to note that, for any node  $(a, n) \in \mathcal{V}$  that is not a sink node, it must hold that  $deg^{out}(a) = n$ , because the number of cases that pass through the node must also continue to another node. Similarly, for any node  $(a, n) \in \mathcal{V}$  that is not a source node, it must hold that  $deg^{in}(a) = n$ . We call these constraints the *process graph rule*.

## 3. Phase I: Discovery of Process Model from Event Logs

The core idea behind the first phase of correlation miner is based on the process graph rule, defined in the previous section, which states that the number of cases

that pass through an activity must be equal to the number of incoming and outgoing cases for that activity. We can count the number of cases that pass through an activity  $a$  by counting the number of occurrences  $|E_a|$  of that activity in the log. Subsequently, we can draw the edges between the activities, in such a way that process graph rule is met. This is a constraint programming problem that can be solved using integer linear programming. At this point we assume that the source and sink nodes are set manually.

As an example, Fig. 2a shows a business process model that can be constructed from the log in Table 2. In this log  $A$  occurs 10 times, while  $D$  occurs 4 times. We manually select  $A$  as a source node. Consequently, it has an in-degree of 0. The number of occurrences of  $A$  corresponds to the out-degree of  $A$ , which is 4 (to  $D$ ) plus 6 (to  $E$ ) equals 10. Similarly, the number of occurrences of  $D$  corresponds to the in-degree and the out-degree of  $D$ . In this way, creating a business process model that meets the process graph rule for all activities, produces Fig. 2a. However, Fig. 2b, 2c, and 2d are also models that meet this rule. While Fig. 2c is the model that matches the original log from Table 1 best, the process graph rule alone is not enough to determine that.

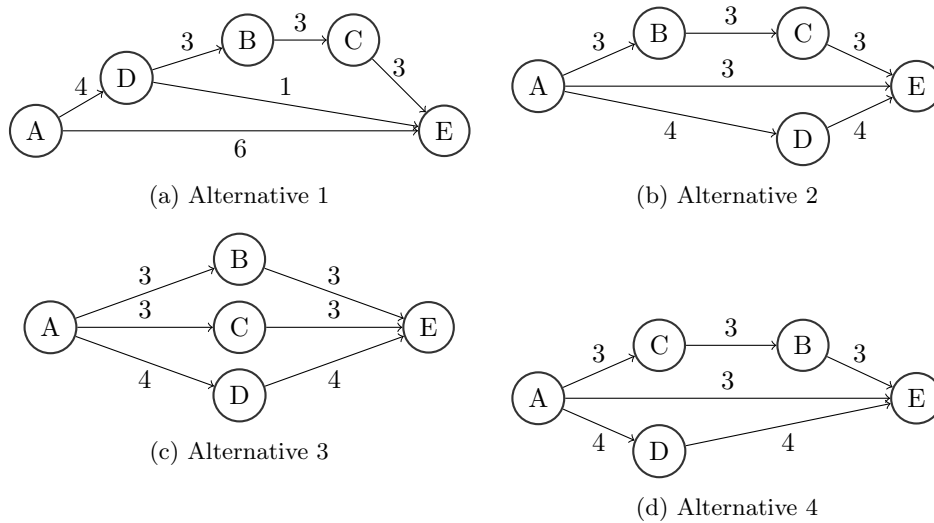


Fig. 2: Feasible process models from the example log

As the example shows, it is often possible to create more than one model that meets the criteria. Therefore, additional measurements from the event log are required in order to select the best model. To this end, the correlation miner creates two matrices, the *Precede/Succeed matrix* and the *Duration matrix*. Each score in Precede/Succeed matrix,  $P/S_{i,j}$ , indicates the fraction of events referring

6 *S. Pourmirza, R. Dijkman & P. Grefen*

to activity  $i$  that have occurred before events referring to activity  $j$ . If this score is high, it is more likely that there is an edge from  $i$  to  $j$ . Each score in Duration matrix,  $D_{i,j}$ , indicates the average time or standard deviation (we will experiment with both alternatives) of the time difference between events referring to activity  $i$  and events referring to activity  $j$ . If this score is low, it is more likely that there is an edge from  $i$  to  $j$ . In the remainder of this section, we explain how these matrices are computed, how they are used to construct the linear programming problem, and ultimately how the business process model is constructed.

### 3.1. Precede/Succeed Matrix

The first step in the first phase of correlation miner is to calculate the Precede/Succeed matrix. This matrix is a square matrix of order  $n$ , with  $n = |A|$ . Let  $i$  represent a row in the matrix as well as an activity from  $A$  and let  $j$  represent a column and an activity. A value  $P/S_{i,j}$  in the matrix represent the fraction of events from  $i$  that occurred before events from  $j$  and is computed as follows.  $\Omega_{i,j} = \{(e, f) | e \in E_i \wedge f \in E_j\}$  is the set of all pairs of occurrences of  $i$  and  $j$ . Furthermore, given events  $e$  and  $f$ , let  $b(e, f)$  be a function that is 1 if  $t_e < t_f$  and 0 otherwise.

$$P/S_{i,j} = \frac{\sum_{(e,f) \in \Omega_{i,j}} b(e, f)}{|\Omega_{i,j}|}$$

Returning to our example (Table 2),  $P/S_{C,E}$  can be computed as follows.  $\Omega_{C,E}$  contains 30 elements since  $|E_C| = 3$  and  $|E_E| = 10$ .  $\sum_{(e,f) \in \Omega_{C,E}} b(e, f) = 17$ , because in 17 cases  $C$  occurred before  $E$ . Therefore,  $P/S_{C,E} \approx 0.57$ .

Similarly, we can calculate the value for each pair of activities in a given event log. Fig. 3 illustrates the matrix that is calculated based on the example log.

For performance reasons, we can remove some pairs to not include them for consideration during the next two steps of the algorithm. This can be done via a threshold filter. The goal of the threshold filter is to remove pairs of activities that have a very low probability of forming an edge. For example, in Fig. 3, we have  $P/S_{A,C} = 0.97$  and  $P/S_{C,A} = 0.03$ . These values indicate that it

$$P/S = \begin{matrix} & A & B & C & D & E \\ \begin{matrix} A \\ B \\ C \\ D \\ E \end{matrix} & \begin{pmatrix} 0.00 & 0.47 & 0.97 & 0.73 & 0.88 \\ 0.53 & 0.00 & 1.00 & 0.67 & 0.90 \\ 0.03 & 0.00 & 0.00 & 0.33 & 0.57 \\ 0.26 & 0.33 & 0.67 & 0.00 & 0.70 \\ 0.12 & 0.10 & 0.43 & 0.30 & 0.00 \end{pmatrix} \end{matrix}$$

Fig. 3: P/S matrix for the example log

is likely that there is an edge from  $A$  to  $C$ , but no edge from  $C$  to  $A$ . The challenge is to select a threshold that removes as many pairs as possible in order to increase performance, but to prevent false positives, because pairs that are removed at this

stage will never become edges in the business process model.

### 3.2. Duration Matrix

The second step in the first phase of correlation miner is to generate the Duration matrix. This matrix is a square matrix of order  $n$ , with  $n = |A|$ . Let  $i$  represent a row in the matrix as well as an activity from  $A$  and let  $j$  represent a column and an activity. A value  $D_{i,j}$  indicates the average time difference or the standard deviation of the time difference between events referring to activity  $i$  and events referring to activity  $j$ .

In order to calculate the value of  $D_{i,j}$  we need to have a *mapping* between elements of  $E_i$  and elements of  $E_j$ . The idea behind this mapping is that it maps events that belong to the same case. We compute the mapping, by relating the events in such a way that the variance in the time difference is as low as possible, because we argue that the time difference between events from the same case have a more constant probability distribution than the time difference between events from different cases. This principle is not likely to yield a perfect mapping, in the sense that events are mapped if and only if they belong to the same case. However, a perfect mapping is not necessary. We just need a mapping that has a time distribution that is close enough to the perfect mapping.

Formally, let  $\Omega_{i,j}$  be defined as in Section 3.1. Then we compute the mapping  $M_{i,j} \subseteq \Omega_{i,j}$  that satisfies the following constraints.

$$\begin{aligned} \forall (e, f) \in M_{i,j} : \exists (e, f') \in M_{i,j} &\Rightarrow f = f' \\ \forall (e, f) \in M_{i,j} : \exists (e', f) \in M_{i,j} &\Rightarrow e = e' \\ \forall (e, f) \in M_{i,j} : t_e < t_f & \end{aligned}$$

These constraints state that each event can be mapped at most once and that each event from  $E_i$  must be mapped to an event from  $E_j$  that occurs at a later point in time. In addition, the mapping must maximize its size (i.e. as many events must be mapped as possible), while minimizing the standard deviation of the time difference between mapped events. With this mapping, we can compute  $D_{i,j}$  as either the mean  $\overline{\Delta M_{i,j}}$  or the standard deviation  $\sigma_{\Delta M_{i,j}}$ .

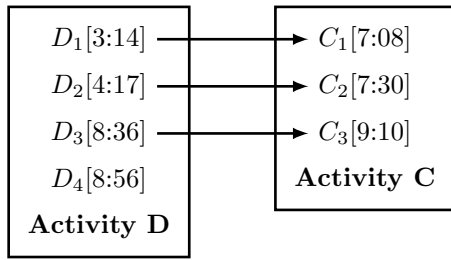
$|M_{i,j}|$  must be maximized

$\sigma_{\Delta M_{i,j}}$  must be minimized, where  $\Delta M_{i,j} = \{t_f - t_e | (e, f) \in M_{i,j}\}$

With this mapping, we can compute  $D_{i,j}$  as either the mean  $\overline{\Delta M_{i,j}}$  or the standard deviation  $\sigma_{\Delta M_{i,j}}$ . Returning to our example (Table 2),  $D_{D,C}$  can be computed as follows. Fig. 4 shows the mapping from events  $E_D$  to events  $E_C$  that



satisfies the constraints above. The average time difference between mapped events is approximately 153 seconds. Similarly, we can calculate the value for each pair of activities in a given event log. Fig. 5 illustrates the matrix that is calculated based on the example log. Note that the time difference in the matrix is slightly different from the time difference that we computed above. This is because we implemented an efficient (greedy) algorithm to compute the mapping, which does not always produce the mapping that satisfies all constraints. In particular, it does not always return the mapping with the lowest standard deviation, but rather the mapping with a local minimum. Although a lower score of  $D_{i,j}$  indicates that it is more likely that there is an edge from  $i$  to  $j$ , a score of 0 means it is not possible to have such an edge, since we assume that events cannot happen at the same time in the event log.


 Fig. 4: Mapping of  $E_D$  to  $E_C$ 

$$D = \begin{matrix} & A & B & C & D & E \\ \begin{matrix} A \\ B \\ C \\ D \\ E \end{matrix} & \begin{pmatrix} 0 & 74 & 109 & 108 & 220 \\ 54 & 0 & 214 & 164 & 69 \\ 6 & 0 & 0 & 87 & 42 \\ 0 & 106 & 150 & 0 & 124 \\ 0 & 102 & 62 & 96 & 0 \end{pmatrix} \end{matrix}$$

Fig. 5: Duration matrix of example log

### 3.3. Business Process Model Construction

The final step in the first phase of correlation miner is to generate the business process model. We first explain how all possible business process models can be generated, taking the process graph rule into account. Then we explain how the best business process model can be selected from all possible business process models, using the Precede/Succeed matrix (Section 3.1) and Duration matrix (Section 3.2).

In order to construct feasible models, we formulate our problem as an *Integer Linear Programming* (ILP) problem. Each possible edge  $(i, j)$  between activities becomes a variable  $x_{ij}$  of our ILP problem. The value that is assigned to each variable,  $x_{ij}$ , indicates the frequency with which an event for activity  $i$  is directly followed by an event from activity  $j$  in a case. Therefore this value can be seen as an edge weight for the edge that connects the node that refers to activity  $i$  to the node that refers to activity  $j$ . Let  $\mathcal{X}$  denote the set that contains the introduced variables. Also, let  $A_s$  be the set of activities that we consider start activities and  $A_e$  be the set of activities that we consider end activities. We can reduce the problem space, by removing all variables:

- $x_{ij}$  that represent unlikely edges, because of a  $P/S_{i,j}$  or  $D_{i,j}$  thresholds;

- $x_{is}, s \in A_s$ , which represent incoming edges to start activities; and
- $x_{ei}, e \in A_e$ , which represent outgoing edges from end activities.

Subsequently, for the resulting set of variables  $\mathcal{X}$ , we can formulate the constraints of our ILP problem as follows. The lower bound for each variable is 0, while the upper bound for each variable is  $\min(|E_i|, |E_j|)$ , because these are the minimum and the maximum number of cases that can flow on the edges. The sum of the number of cases on the incoming edges of an activity  $a$ , must be equal to the number of times an event  $E_a$  occurs for that activity. Similarly, the sum of the number of cases on the outgoing edges of an activity  $a$ , must be equal to the number of times an event  $E_a$  occurs for that activity. Consequently, the constraints become:

$$\begin{aligned}
 x_{ij} &\geq 0 && \text{for each } a \notin A_s : \sum_{x_{ia} \in \mathcal{X}} x_{ia} = |E_a| \\
 x_{ij} &\leq \min(|E_i|, |E_j|) && \text{for each } a \notin A_e : \sum_{x_{ai} \in \mathcal{X}} x_{ai} = |E_a|
 \end{aligned}$$

Returning to the example of the log from Table 2, we can formulate the following constraints.  $x_{AE} \geq 0$  and  $x_{AE} \leq 10$ . Similarly,  $x_{AB} \geq 0$  and  $x_{AB} \leq 3$ . Also,  $x_{AE} + x_{BE} + x_{CE} + x_{DE} = 10$  and  $x_{AB} + x_{AC} + x_{AD} + x_{AE} = 10$ . In addition,  $x_{CB} + x_{CD} + x_{CE} = 3$ , and  $x_{AC} + x_{BC} + x_{DC} = 3$ . Note that, since  $x_{CA}$  represents an incoming edge to a start node  $A$  and  $x_{EC}$  represents an outgoing edge from an end node  $E$ , these two variables have been removed and consequently are ignored in these constraints.

By applying these constraints, we can construct all feasible business process models. Now in the second part of this section, our goal is to select the model, from all feasible ones, that best represents the event log behavior. For this purpose, we employ matrices that we computed in the first two steps in order to formulate an *objective function* to eventually select the *optimal* model out of all feasible ones, using ILP principles.

The matrices from step 1 and step 2 provide some evidence to suggest that if activity pair  $(i, j)$  contains a high value for  $P/S_{i,j}$  and a low value for  $D_{i,j}$ , these two activities may have a higher chance to form an edge. Based on this statement, our objective is to select the model that has a higher cumulative value of  $P/S$ s and a lower cumulative value of  $D$ s for the edges. To achieve this objective, we define a *coefficient* for each variable  $x_{ij}$ , in our ILP in order to consider these values. Since we are interested in a high value of  $P/S_{i,j}$  and a low value of  $D_{i,j}$ , we define an *edge ratio* as  $\frac{D_{i,j}}{P/S_{i,j}}$  for each variable  $x_{ij}$  that must be minimized. Note that, we need to keep the sum of edge ratios for each edge independent of the assigned value of the referring variables and count the edge ratio if there is an edge, irrespective of the weight of that edge. Thus, ideally, each  $x_{ij}$  is divided by its value and set to 0 if  $x_{ij} = 0$  in order to eliminate their influence. However, this is not possible in an ILP. Therefore, in order to reduce this influence we divide each variable by

10 *S. Pourmirza, R. Dijkman & P. Grefen*

its upper bound. In conclusion, we formulate the objective function for our ILP as follow:

$$C_{ij} = \frac{D_{i,j}}{P/S_{i,j}} \cdot \frac{1}{\min(|E_i|, |E_j|)}$$

$$\text{minimize } \sum_{x_{ij} \in \mathcal{X}} C_{ij} \cdot x_{ij}$$

Returning to the example log (Table 2), by considering the constraints, multiple models can be constructed as shown on Fig. 2. Now, we can calculate the objective function for each of these business process models and select the one that has a lowest value as an output of our algorithm. The value for the business process model from Fig. 2c is the lowest. Therefore, we select the business process model of Fig. 2c as the final result for the first phase of correlation miner.

The algorithm that we described above may return a model that contains cycles. However, since the correlation miner only deals with acyclic processes, these cycles should be removed. Therefore, it employs Johnson's algorithm<sup>13</sup> to capture all edges that are involved in cycles in the model. It then calculates the edge ratio for these edges and removes the edge with highest edge ratio. The algorithm then reruns the third step repetitively until a model without cycles is captured. Note that, as the main loop of the algorithm is ILP solving, the time complexity is polynomial.

#### 4. Evaluation of Phase I

This section presents the evaluation of the first phase of correlation miner. It first presents the setup of the evaluation in Section 4.1. Then it presents two evaluations, one using synthetic event logs (Section 4.2) and one using a modified version of a real-world event log (Section 4.3).

##### 4.1. Evaluation Setup for Phase I

In order to conduct our evaluation, we have implemented the first phase of correlation miner as a Java application. Since the algorithm uses the Gurobi ILP Solver<sup>14</sup>, which is commercial software, one needs to install this software on his machine with a proper license. For the interested reader we provide on-line access to the algorithm<sup>a</sup>, which can only be used for the academic purposes. Moreover, we are currently investigating the licensing issue in order to implement the algorithm as a ProM<sup>15</sup> plug-in.

Fig. 6 depicts the procedure for evaluating our algorithm. For the evaluation with synthetic logs, we first create synthetic business process models with particular properties, to investigate the effect that these properties have on our algorithm.

<sup>a</sup><http://is.ieis.tue.nl/research/correlation>

This will be explained in detail in Section 4.2. Subsequently (step 1a), we generate synthetic logs for these models, using the BIMP simulator<sup>b</sup>. For the evaluation with real-world logs, we take a log from practice. Subsequently (step 1b), we generate a business process model for that log. This will be explained in detail in Section 4.3. Now that we have both a log and a model, we remove the case identifiers from the log (step 2) to generate a log that can be used for correlation mining and (step 3) we apply the first phase of correlation miner. Finally (step 4), we assess the quality of the mined model.

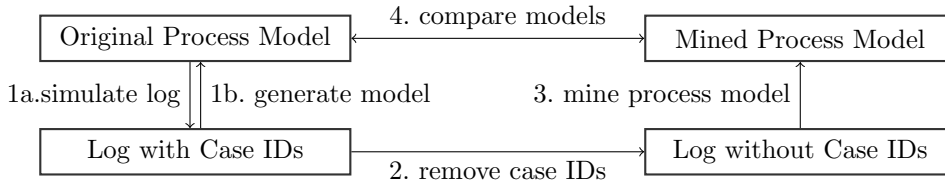


Fig. 6: Evaluation procedure of mined process models

In order to measure the quality of the mined models, we use *precision* and *recall*. Precision measures the fraction of edges in the mined model that are correct, i.e. are also in the original business process model. Recall measures the fraction of correct edges that have been found. These measures are defined as follows. Let  $TP$  be the set of edges that exist in the mined model and also in the original;  $FN$  be the set of edges that do not exist in the mined model but do exist in the original model; and  $FP$  be the set of edges that exist in the mined model but do not exist in the original model. Then precision and recall are defined as:

$$precision = \frac{TP}{TP + FP} \quad recall = \frac{TP}{TP + FN}$$

Returning to Fig. 2 and assuming that Fig. 2b is the mined model and Fig. 2c the original model, we can calculate precision and recall as follows.  $TP$  for these two models is 4, including  $AB$ ,  $AD$ ,  $CE$  and  $DE$ ;  $FN$  is 2, including  $AC$  and  $BE$ ; and finally  $FP$  is again 2, including  $BC$  and  $AE$ . Consequently, precision and recall for Fig. 2b are approximately 0.67.

In related work fitness and appropriateness have been introduced<sup>4</sup> as measures to evaluate the quality of a mined business process model. However, these measures evaluate the quality of a business process model as it is compared to a log. We evaluate the quality of the business process model as it is compared to another business process model; the business process model that should have been returned. By doing so, we can get more meaningful results. However, clearly this is only

<sup>b</sup><http://bimp.cs.ut.ee/>

possible when such a business process model indeed exists. In other situations fitness and appropriateness should be used.

#### 4.2. *Synthetic Event Logs for Phase I*

We evaluated the first phase of correlation miner using over 250 synthetic event logs with different properties. We varied properties that we assumed would have an effect on the quality of the mined business process model. Specifically, we generated logs:

- using both structured and unstructured business process models (according to the criteria in <sup>16</sup>) and using business process models with different numbers of branches, because these factors influence model complexity and we assume that more complex models are more difficult to mine;
- that contained different numbers of cases, because a higher number of cases provides more data to mine from and should therefore produce more accurate models;
- with different inter-arrival times and activity durations, because time properties play an important role in our algorithm and should therefore have an impact on the quality of the result.

In total we used 6 business process models to generate our logs. For reasons of space, we do not present all results, but rather the results for the model that produced the best results (Fig. 7a) and the model that produced the worst results (Fig. 7b). These models primarily vary with respect to whether they are structured or unstructured. The maximum number of branches does not vary much, because our evaluation showed that the number of branches did not play a substantial role in the quality of the results.

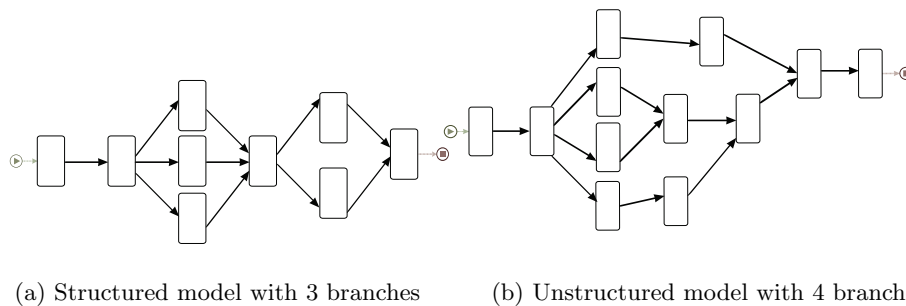


Fig. 7: Evaluation business process models

Table 3 presents the results for mining logs that were generated from the model from Fig. 7a. We generated different logs from this model, in such a way that the duration of activities and the inter-arrival times, were either:

- distributed uniformly with a median that was selected randomly from the interval  $(1, 100]$ ;
- distributed normally with a mean that was selected randomly from the interval  $(1, 100]$  and a standard deviation that was selected randomly from the interval of  $\frac{1}{5}$  to  $\frac{1}{7}$  of the selected mean;
- distributed exponentially with a mean that was selected randomly from the interval  $(1, 100]$ .

For each of these three distributions, we produced an event log with 200, 2,000, and 10,000 cases, which is the maximum number of cases that can be produced by the BIMP Simulator. In two experiments, the correlation miner was not able to generate a business process model, because after removing some variables (e.g. by removing cycles, or P/S threshold filter), it was unable to solve the model with ILP, as the model became infeasible with regard to its constraints. The results in this table show that the correlation miner can discover a simple structured model perfectly, provided that there are sufficiently many cases in the log.

Table 3: Results of using the first phase of correlation miner for Fig. 7a

Distribution	Uniform			Exponential			Normal		
# of cases	200	2000	10000	200	2000	10000	200	2000	10000
Model Recall	1.00	1.00	1.00	n/a	1.00	1.00	n/a	1.00	1.00
Model Precision	1.00	1.00	1.00	n/a	1.00	1.00	n/a	1.00	1.00

Table 4 presents an overview of the results for mining logs that were generated from the model from Fig. 7b. Experiments on these logs have been conducted with activity duration distributions that were selected in the same manner as for the model from Fig. 7a, but with inter-arrival times that were either selected in the same manner or from an interval with a longer duration  $((200, 300])$ . We also experimented with the ranges  $(1, 10]$  and  $[10, 10]$ . However, these did not lead to substantially different results. Therefore, due to space restrictions, we do not publish those results here. For each of these combinations of probability distributions, we produced an event log with 100, 1,000, and 10,000 cases.

Under most conditions, the results for this more complex unstructured model are worse than for the structured model. If a structured model contains a task with  $n$  outgoing flows, there is for sure another task with  $n$  incoming flows; however, if an unstructured model contains a task with  $n$  outgoing flows, it may not have any task with  $n$  incoming flows. Therefore, solving the structured model is easier for our ILP. The other interesting finding to emerge from Table 4 is that if the inter-arrival time between cases is higher than the service time of activities, this leads to substantially better results. Comparing the two sets of results in this table,

shows that in the experiments with higher inter-arrival time the results for recall and precision noticeably increased on average by 27% and 42% respectively. This result can be explained, because if the inter-arrival time is the same as the service time, cases are more ‘intertwined’ in the log, such that the correlation miner has more difficulties telling them apart based on their timing properties (which other mining algorithms can do based on case identifiers).

Table 4: Results of using the first phase of correlation miner for Fig. 7b

Inter-arrival time	[1-100]					
Distribution	Uniform			Exponential		
# of cases	100	1000	10000	100	1000	10000
Model Recall	0.57	0.79	0.64	0.57	0.57	0.71
Model Precision	0.42	0.65	0.47	0.42	0.42	0.56
Inter-arrival time	[200-300]					
Model Recall	0.75	0.93	1.00	0.93	0.71	1.00
Model Precision	0.86	0.87	1.00	0.81	0.59	1.00

Also, as expected, and in line with the findings from Table 3, a higher number of cases in the log leads to better models in most cases. However, in experiments in which activities were distributed exponentially, our algorithm did not find a significant difference between the event logs that originally contained 100 and 1000 cases.

Hence, based on these results, we can conclude that the correlation miner is applicable to discover process models from event logs without case identifiers, especially for structured models and when the inter-arrival time between cases is higher than the service time of activities and the number of cases in the log is sufficiently high.

### 4.3. Real-World Event Log for Phase I

We also evaluated the algorithm based on a real-world event log. To this end, we used a modified version of the log from the BPI Challenge of 2012<sup>17</sup>. First, we mined a business process model from this event log using Disco<sup>9</sup>. We then removed all loops in order to make it acyclic, because that our algorithm only deals with acyclic models at this stage. For the same reason, we then removed the cases from the log in which loops appeared. The resulting version of this real-world event log contained 13 activities, 70,425 events, and 11,647 cases.

In accordance with the evaluation procedure described in Section 4.1, we then

removed all case identifiers from the log, discovered process models from the log using the correlation miner, and subsequently, evaluated the quality of the discovered model. The results obtained from this experiment show a precision of 85% and a recall of 63%. The lower result for recall can be explained by the fact that our algorithm seeks to find edges with higher frequencies, while we used the Disco model that included all possible edges that could be mined from the log, not just the ones with the higher frequency. Normally speaking, one would not be interested in the Disco model that contains all edges, because this model also includes edges that represent exceptional situations. Indeed if we remove edges with lower frequency (i.e.,  $< 50$ ) from the original model our recall indeed increases to 68% as precision decreases to 79%.

Based on these results, we claim that the correlation miner is also applicable to mine real-world event logs.

## 5. Phase II: Reconstruction of Cases from Event Logs

The first phase of correlation miner (cf. Section 3), allows us to construct a process model from an event log without case identifiers. However, it does not allow us to reconstruct the cases themselves; it does not reveal which events belong to the same case. This section fills that gap, by presenting the second phase of correlation miner aiming to accurately determine which events belong to the same case.

The algorithm proposed in this section works by interpreting the log as a set of binary relations between events. A relation from event  $e$  to event  $f$  represents that event  $e$  is directly followed by event  $f$  in the same case. We also say that  $e$  caused  $f$ . For example, considering the event log shown in Table 1, Case 1 contains:

- a binary relation between the event referring to activity  $A$  that happened at 00:20 and the event referring to activity  $B$  that happened at 02:04, depicted by  $(A_1[00:20], B_1[02:04])$ , and
- a binary relation between the event referring to activity  $B$  that happened at 02:04 and the event referring to activity  $E$  that happened at 02:32, depicted by  $(B_1[02:04], E_1[02:32])$ .

Therefore, Case 1 can be denoted as  $\{(A_1[00:20], B_1[02:04]), (B_1[02:04], E_1[02:32])\}$ . Alternatively, the relations can be represented in a graph. For example, the graph that represents Case 1 is illustrated in Figure 8.

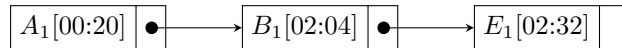


Fig. 8: Graph representation case 1 from the example event log

Against this background, the goal of the second phase of correlation miner is to find the binary relations that determine the events that directly follow each other in



a case. We achieve that goal by casting the problem into a Quadratic Programming (QP) problem, in which we define constraints that must hold for the binary relations, and an optimization function that defines what the ‘best’ set of binary relation is.

In the remainder of this section we describe how the mentioned constraints and objective function are formulated, in Section 5.1 and 5.2 respectively.

### 5.1. Constraints

We present a number of constraints that must hold for causal relations between events. These constraints help us reduce the space of possible relations that must be explored. In particular, the constraints that must hold are that:

- two events can only be causally related, if their activities are also related in the mined process model; and
- an event can be directly preceded or directly followed by at most one other event.

Considering a discovered process model that includes an edge from a node referring to activity  $i$  to a node referring to activity  $j$ . There is a possibility that  $e_i \in E_i$  was directly followed by  $e_j \in E_j$ , and therefore, the referring binary relation  $(e_i, e_j)$  must exist in one of the execution cases. In general, each binary relation, e.g.,  $(e_x, e_y)$ , introduces a binary variable,  $r_{e_x e_y}$ , in our QP problem, which indicates whether the referring binary relation exists in a case or not. Clearly, we can reduce the problem space, by removing all variables  $r_{e_x e_y}$  that represents a binary relation in which  $e_x$  has occurred after  $e_y$ . Let  $G = (\mathcal{V}, \mathcal{E}, \omega)$  be a discovered process model and  $\mathcal{R}$  denote the set that contains all the introduced variables. We formulate the first set of constraints of our QP problem as follows.

$$\text{for each edge } (x, y) \in \mathcal{E}: \quad \sum_{r_{e_x e_y} \in \mathcal{R}} r_{e_x e_y} = \omega(x, y)$$

Returning to the example of the log from Table 2 and the discovered process model shown in Figure 2c, we can, for example, conclude that there must be three binary relations between the events referring to activity  $C$  and events referring to activity  $E$  out of 17 possible relations as shown in Figure 9. Therefore, we can formulate the following constraint for the edge between the node referring to activity  $C$  and the node referring to activity  $E$ .

$$r_{C_1 E_4} + \dots + r_{C_1 E_{10}} + r_{C_2 E_5} + \dots + r_{C_2 E_{10}} + r_{C_3 E_8} + \dots + r_{C_3 E_{10}} = \omega(C, E) = 3$$

Obviously, if the value of  $r_{C_1 E_4}$  is 1 (i.e.,  $C_1$  and  $E_4$  belonged to the same case), no other binary relation can contain  $C_1$  as a source event and  $E_4$  as a target event. Therefore, to apply this restriction, we introduce two additional sets of binary variables for each event to specify whether the referring event has already been selected as (i) a source event, and as a (ii) target event. Let  $\mathcal{S}_{e_a e_x}$  denote the set that contains all the binary variables referring the usage of  $e_a$  as a source event, and  $\mathcal{T}_{e_x e_a}$  denote the set that includes all the binary variables referring the usage

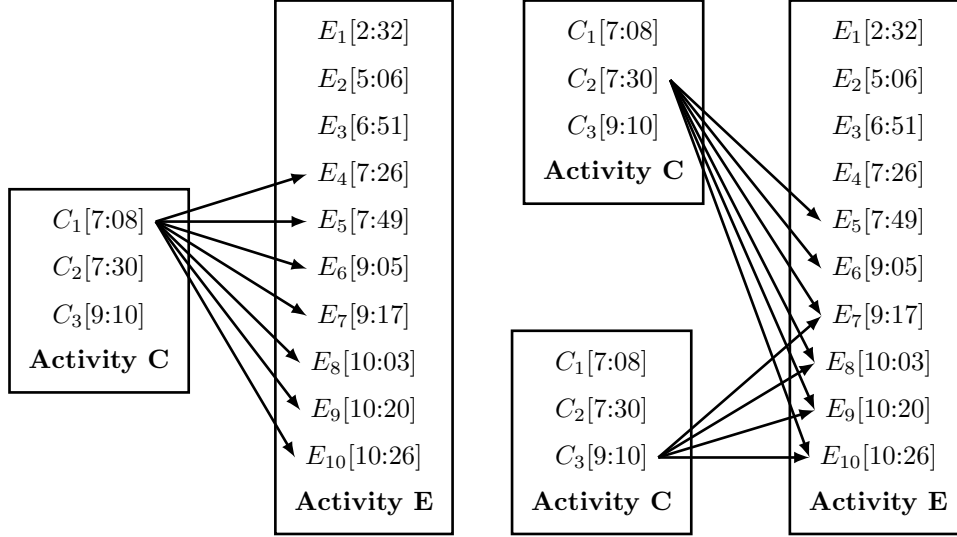


Fig. 9: All possible binary relations between events referring to activity  $C$  and events referring to activity  $E$

of  $e_a$  as a target event. In addition, let  $E_s$  be the set of events referring to the start activities and  $E_e$  be the set of activities referring to the end activities. We can reduce the problem space, by removing all variables

- $t_{e_x e_s} \in \mathcal{T}_{e_x e_s}$  such that  $e_s \in E_s$ , which represent the usage of events referring to start activities as target events; and
- $s_{e_e e_x} \in \mathcal{S}_{e_e e_x}$  such that  $e_e \in E_e$ , which represent the usage of events referring to end activities as source events.

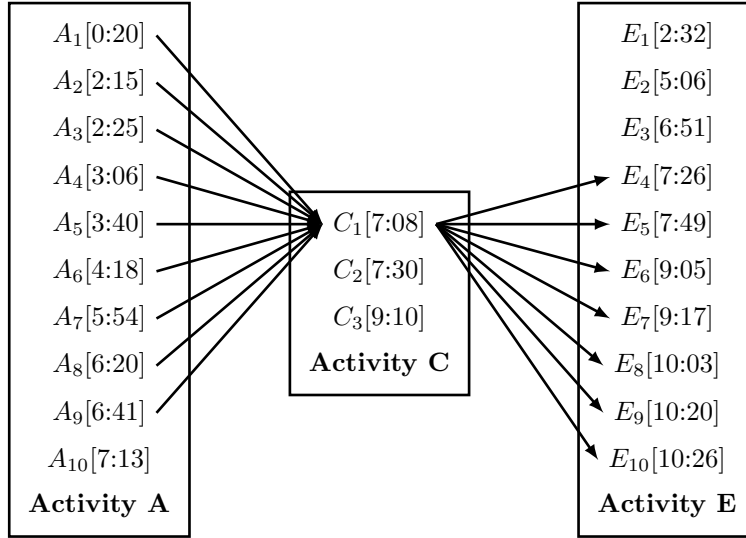
Subsequently, for the resulting sets of binary variables, we can formulate the constraints of our QP problem as follows.

$$\begin{aligned} \text{for each event } e_a \in E_a \wedge E_a \neq E_s: & \quad \sum_{t_{e_x e_a} \in \mathcal{T}_{e_x e_a}} t_{e_x e_a} = 1 \\ \text{for each event } e_a \in E_a \wedge E_a \neq E_e: & \quad \sum_{s_{e_a e_x} \in \mathcal{S}_{e_a e_x}} s_{e_a e_x} = 1 \end{aligned}$$

Returning to the example of the log from Table 2 and the discovered process model shown in Figure 2c, we can, for example, conclude that  $C_1$  must be a target event for only one of the 9 possible relations as shown on the left side of Figure 10, and similarly, it requires to be a source event for only one of the 7 possible relations as shown on the right side of Figure 10. Therefore, we can formulate the following constraints for this event.

$$t_{A_1 C_1} + t_{A_2 C_1} + t_{A_3 C_1} + t_{A_4 C_1} + t_{A_5 C_1} + t_{A_6 C_1} + t_{A_7 C_1} + t_{A_8 C_1} + t_{A_9 C_1} = 1$$

$$s_{C_1E_4} + s_{C_1E_5} + s_{C_1E_6} + s_{C_1E_7} + s_{C_1E_8} + s_{C_1E_9} + s_{C_1E_{10}} = 1$$


 Fig. 10: All possible situation for the first event of  $C$ 

Accordingly, the value of  $r_{C_1E_4}$  indicates the values of  $s_{C_1E_4}$  and  $t_{C_1E_4}$ . In other words, if  $C_1$  and  $E_4$  belong to the same case, there is only one binary relation in which  $C_1$  acts as a source event that is  $s_{C_1E_4}$ , and only one binary relation in which  $E_4$  acts as a target event that is  $t_{C_1E_4}$ . Conversely, if these two events belong to different cases, the values of  $s_{C_1E_4}$  and  $t_{C_1E_4}$  are 0. In general, we can state that the value for variable  $r_{e_xe_y} \in \mathcal{R}$  will be 1, if and only if there exists

- a variable  $s_{e_xe_y} \in \mathcal{S}_{e_xe_y}$  such that  $s_{e_xe_y} = 1$ , and
- a variable  $t_{e_xe_y} \in \mathcal{T}_{e_xe_y}$ , such that  $t_{e_xe_y} = 1$ .

To ensure that the value of a binary relation is the same as the values of referring source and target usage variable, we formulate the following constraint in our QP problem.

$$\text{for each two events } e_a \in E_a \text{ and } e_b \in E_b: \quad r_{e_ae_b} = s_{e_ae_b} = t_{e_ae_b}$$

Finally, note that, these constraints are formulated based on results of the first phase of correlation miner. Consequently, if a discovered process model does not precisely depict the behavior of an event log, these constraints may result into an infeasible QP problem. In order to mitigate this threat, we allow the constraints to be violated according to predefined penalty costs. For instance, the constraints referring to captured process model (e.g., first set of constraints) have the lowest penalty costs (i.e., they can be easily violated). This decision decision can be

explained by the fact that, in practice, there is no ground truth process model for an event log since event logs only contain a fraction of possible cases (i.e., many other cases can possibly be occurred) and also event logs cannot show impossible cases (i.e., cases that cannot be occurred). In contrast with this penalty cost, the constraints referring to usage of an event have the highest penalty cost since, as noted before, events can be directly preceded or directly followed by at most one other event.

## 5.2. Objective Function

By applying the constraints in Section 5.1, we can construct all feasible sets of ‘directly follows’ relations between events from the same case. However, we aim to select the best set. For this purpose, we further employ the time difference between the source and target event of each binary relation in order to formulate an objective function to pick the optimal set of relations.

More specifically, this objective function attempts to form a subset of all possible binary relations in such a way that the variation in the time difference among the selected binary relations is as low as possible. This objective function is designed to comply with the assumption that time differences between binary relations that belong to the same cases are more likely to follow a certain probability distributions, than time differences between binary relations from different cases.

The most obvious coefficient that can be used to achieve the mentioned objective is the actual time difference between the source and destination event for each binary relation. However, to precisely define a coefficient for the binary relations, we further adopt the idea of correlation coefficient suggested by Bradley Efron<sup>18</sup>, which is based on the variance and standard error of an unknown distribution for an arbitrary sample. Based on this idea, we define error as the squared difference between the observed data and the expected value. Although in general the mean of a sample can be used as an (expected) parameter that minimizes the total error, in the context of this study, we attempt to minimize the total error by minimizing the squared difference between the source and target event for each binary relation. In other words, we expect that the variation between the source and target event of binary relations is as low as possible. Consequently, we formulate the objective function for our QP problem as follows. Let  $O_{e_a}$  denote the timestamp at which event  $e_a$  has occurred.

$$\text{minimize } \sum_{r_{e_i e_j} \in \mathcal{R}} (O_{e_j} - O_{e_i})^2 \left( \frac{r_{e_i e_j} (O_{e_j} \cdot t_{e_i e_j} - O_{e_i} \cdot s_{e_i e_j})}{O_{e_j} - O_{e_i}} \right)$$

By way of illustration, let us consider the simple example given below. Table 5 shows the original event log, Table 6 exemplifies this event log without case identifiers and Figure 2 presents the discovered process model using the the first phase of correlation miner.

This section set out to determine the cases in the original log (Table 5) by

20 *S. Pourmirza, R. Dijkman & P. Grefen*

Case	Activity	Timestamp
1	X	02:00
1	Y	04:00
1	Z	08:00
2	X	03:00
2	Y	05:00
2	Z	09:00

Table 5: Simple event log with cases identifiers

Activity	Timestamp
X	02:00
X	03:00
Y	04:00
Y	05:00
Z	08:00
Z	09:00

Table 6: Simple event log without case identifiers

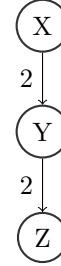


Fig. 11: Mined process model

employing the event log without case identifiers (Table 6) and discovered process model (Figure 11). Looking closely at the events in the log without case identifier, there are three possible sets of binary relations that are applicable for the discovered model. For brevity, we use the shortened form of the timestamps of the mentioned events in the square brackets. These three possible sets are:

- $\{(X_1[2], Y_1[4]), (Y_1[4], Z_1[8]), X_2[3], Y_2[5]), (Y_2[5], Z_2[9])\}$ , which results in two execution cases:  $\langle X_1, Y_1, Z_1 \rangle$  and  $\langle X_2, Y_2, Z_2 \rangle$ ;
- $\{(X_1[2], Y_2[5]), (Y_2[5], Z_1[8]), X_2[3], Y_1[4]), (Y_1[4], Z_2[9])\}$ , which results in two execution cases:  $\langle X_1, Y_2, Z_1 \rangle$  and  $\langle X_2, Y_1, Z_2 \rangle$ ; and
- $\{(X_1[2], Y_2[5]), (Y_2[5], Z_2[9]), X_2[3], Y_1[4]), (Y_1[4], Z_1[8])\}$ , which results in two execution cases:  $\langle X_1, Y_2, Z_2 \rangle$  and  $\langle X_2, Y_1, Z_1 \rangle$ .

Having identified these three sets of binary relations, we can calculate the objective function for each of them to select the one that has the lowest value as the best candidate. The value of the objective function for the first set is 40, for the second set is 44, and for the last set is 42. Therefore, the first set of execution cases has the lowest value, and thus, this set of binary relations will be selected to produce the event log with cases. As confirmed by the original log (Table 5), the selected set of binary relations best reveals the actual ‘directly follows’ relations among the events within the event log without case identifiers (Table 6). In the same way, we can compute the best set of binary relations among the events in Table 2 by employing the discovered process model shown in Figure 2c. This calculation results in exactly the same event log as shown in Figure 1.

## 6. Evaluation of Phase II

This section presents the evaluation of the second phase of correlation miner. First, the evaluation setup is discussed in Section 6.1. Then two types of evaluations are presented, one using synthetic event logs (Section 6.2) and one using a real-world event log (Section 6.3).

### 6.1. Evaluation Setup for Phase II

In order to conduct the evaluation of the second phase of the correlation miner, we have further extended the online version of correlation miner by implementing the algorithm proposed in previous section. The procedure for evaluating the new algorithm is shown in Fig. 12.

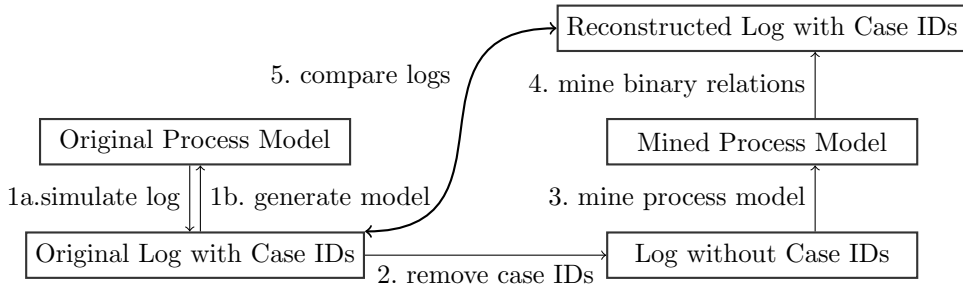


Fig. 12: Evaluation procedure of reconstructed event logs with cases

The first three steps in this procedure are similar to the first three steps for the evaluation procedure of previous algorithm as shown in Fig. 6. However, this evaluation procedure goes on to reconstruct a log with case identifiers from the mined process model in step 4. Subsequently, we compare cases within the reconstructed event log with cases within the original event log in step 5.

In order to measure the quality of the reconstructed event log, we adapt the definition of precision and recall to be used for comparing two event logs. In this context, precision measures the fraction of binary relations in the reconstructed log that are correct, i.e., events that also directly follow each other in a case in the original event log with case identifiers, and recall measures the fraction of correct binary relations that have been found. We further adapt the definition of  $TP$ ,  $FN$  and  $FP$  to be applicable for comparing two event logs. Let  $TP$  be the set of binary relations that exist in the reconstructed event log and also in the original event log;  $FN$  be the set of binary relations that do not exist in the reconstructed log but do exist in the original log; and  $FP$  be the set of binary relations that exist in the reconstructed log but do not exist in the original log. Then, precision and recall can be calculated as explained in Section 4. The main reason that we did not utilize fitness and appropriateness<sup>4</sup> as measures to evaluate the quality of a reconstructed event log is that these measures evaluate the quality of a mined business process model on the basis of an original event log, while we evaluate the quality of reconstructed event log on the basis of that original event log.

Returning to the example given in Table 6 and assuming that the original event log consists of  $\{ \langle X_1, Y_1, Z_1 \rangle, \langle X_2, Y_2, Z_2 \rangle \}$  and the reconstructed event log consists of  $\{ \langle X_1, Y_2, Z_2 \rangle, \langle X_2, Y_1, Z_1 \rangle \}$ , we can calculate precision and recall

22 *S. Pourmirza, R. Dijkman & P. Grefen*

as follows.  $TP$  for these two event logs is 2, including  $(Y_1, Z_1)$  and  $(Y_2, Z_2)$ ;  $FN$  is also 2, including  $(X_1, Y_1)$  and  $(X_2, Y_2)$ ; and finally  $FP$  is again 2, including  $(X_1, Y_2)$  and  $(X_2, Y_1)$ . Consequently, both precision and recall for this example are 0.5.

## 6.2. Synthetic Event Logs for Phase II

We evaluated the proposed technique using the synthetic event logs that were introduced in Section 4.2. However, we decreased the number of cases in these event logs, because at this stage the new algorithm has a high time and space complexity, which makes it impractical for use with event logs with a large number of events. In future work, we will work on reducing the time complexity of the algorithm.

We present the results for the events logs that produced the best results and the event logs that produced the worst results.

Distribution	Uniform			Normal			Exponential		
	# of cases	25	50	75	25	50	75	25	50
Process Model Recall	0.92	1.00	1.00	0.92	1.00	1.00	0.75	1.00	1.00
Log Recall	0.89	1.00	1.00	0.83	1.00	1.00	n/a	1.00	1.00
Process Model Precision	0.85	1.00	1.00	0.85	1.00	1.00	0.75	1.00	1.00
Log Precision	0.92	1.00	1.00	0.86	1.00	1.00	n/a	1.00	1.00

Table 7: Results of using the Correlation Miner for Fig. 7a

Table 7 presents the results for mining logs that were generated from the model from Fig. 7a. We generated different logs from this model by varying different properties, as explained in Section 4.2. The table both shows the precision and recall of the mined process model and the precision and recall of the (binary relations in) reconstructed event log.

The table illustrates that precision and recall for both the mined process model and the reconstructed event log are high. As shown in Section 4, since this model is structured (according to the criteria in <sup>16</sup>) the correlation miner is generally able to generate a business process model with a very high value for recall and precision. Since the mined process model is input for mining the relations, this helps to correctly reconstruct the event log.

The ‘n/a’ values in the table are caused by the fact that the correlation miner was not able to complete the second phase for the exponentially distributed experiments with 25 cases. The reason for that was that the generated process model was less accurate, and that the event log consists of very few cases. However, in general, we conclude that the

Table 8 presents an overview of the results for mining logs that were generated

from the model shown in Fig. 7b. The precision and recall for these models are the worst that we obtained. However, generally, they are still quite good at around 0.8 or even 0.9 precision and recall. Again, there seems to be a strong relation between the quality of the mined process model and the quality of the reconstructed event log.

Distribution	Uniform Distribution								
# of cases	25			50			75		
Random Number	[1-10]	[1-100]	[10]	[1-10]	[1-100]	[10]	[1-10]	[1-100]	[10]
Model Recall	0.79	0.86	0.86	0.86	0.86	0.86	0.93	0.93	0.93
Log Recall	0.73	0.84	0.57	0.58	0.59	0.59	0.83	0.82	0.82
Model Precision	0.69	0.80	0.75	0.80	0.75	0.75	0.87	0.87	0.87
Log Precision	0.78	0.90	0.62	0.66	0.63	0.63	0.84	0.84	0.84
Distribution	Exponential Distribution								
Model Recall	0.93	0.86	1.00	0.93	0.93	1.00	0.93	0.93	0.86
Log Recall	0.78	0.79	1.00	0.81	0.80	1.00	0.81	0.82	0.70
Model Precision	0.87	0.75	1.00	0.87	0.87	1.00	0.87	0.87	0.86
Log Precision	0.81	0.80	1.00	0.85	0.83	1.00	0.83	0.83	0.72

Table 8: Results of Using the Correlation Miner for Fig. 7b

The results shown in this table have been produced by using the experiments in which the inter-arrival time between the cases have a longer duration than service times of activities within an event log. We only present the results for this set of experiments, because the correlation miner can only mine unstructured process models from event logs that contain few cases, if the inter-arrival time between these cases is sufficiently high (i.e., cases are less intertwined in the event logs). As we discussed in Section 4, results for more complex and unstructured process models are generally worse than structured ones, and consequently, we can infer that the quality of reconstructed event logs from unstructured models must be substantially lower than the quality of reconstructed event logs from structured ones. In Table 8, we can see that the results for event logs in which activities and inter-arrival time between cases has exponential distributions are slightly better (9% on average) than event logs in which activities and inter-arrival time between cases have uniform distributions. Finally, the data shown in this table suggest that using different ranges in our experiments did not cause substantially different results since, on average, the difference between the ranges leads to less than 1% differences



among the quality of mined process models and approximately 4% differences among the quality of reconstructed event logs. Table 9 summarizes the data presented in Table 8 by presenting the average value for each quality criterion.

Table 9: Summary of data presented in Table 8

Model Recall	Recons. Log Recall	Model Precision	Recons. Log Precision
0.89	0.76	0.82	0.78

Interestingly, what stands out in the presented results in Table 7 and Table 8, as well as in all the other experiments that we conducted is the value of precision and recall for the reconstructed event log, when the value of precision and recall for a mined process model is 1. During our experiments, we never observed a situation where the value of the precision and recall for mined process models is 1 while the value of these two quality criteria for its related reconstructed event logs is lower than 1.

Hence, based on these results, we can conclude that the correlation miner is applicable to mine logs without case identifiers, especially for structured models and when the inter-arrival time between cases is higher than the service time of activities and the number of cases in the log is high enough.

### 6.3. *Real-World Event Log for Phase II*

We also evaluated the log reconstruction algorithm based on the real-world event log introduced in Section 4.3. However, to make this log usable for the algorithm we reduced the size of the log by randomly selecting 90 cases. To factor out effects of the random selection, we ran the experiment four times on four random selections. The results obtained from these experiments are illustrated in Table 10.

Table 10: Results of using the correlation miner for a real-world event log

Run	Model Recall	Log Recall	Model Precision	Log Precision
1	0.62	0.74	0.72	0.78
2	0.57	0.66	0.60	0.69
3	0.59	0.63	0.56	0.70
4	0.55	0.59	0.63	0.65
Average	0.58	0.65	0.63	0.70

It is apparent from this table that recall and precision for the model reconstruction phase of the correlation miner are better than for its process

discovery phase. In other words, although the quality of the mined process models is low – which is not surprising as we only selected 90 cases, which is low for reliably mining process models –, the second phase of correlation miner was able to find a reasonable number of causal relation in the event logs.

As we discussed in 4.3, the value for recall of real-world event log is generally lower due to exceptional cases that have been recorded in event logs. In order to minimize the effect of these exceptions, we shrank the original real-world event by removing the edges that have a low frequency (i.e.,  $< 50$ ), and subsequently, we generated four event logs with randomly selected cases. These event logs resulted in an increase of both reconstructed event logs' recall and precision, on average, to 71% and 74% respectively. The detailed results of these experiments are illustrated in Table 11. Consequently, based on these results, we claim that the event log reconstruction phase of the correlation miner is also applicable to mine real-world event logs.

Table 11: Results of using the correlation miner for a real-world event log after removing less frequent edges

Run	Model Recall	Log Recall	Model Precision	Log Precision
1	0.60	0.74	0.71	0.81
2	0.70	0.74	0.78	0.76
3	0.63	0.69	0.67	0.72
4	0.68	0.68	0.72	0.71
Average	0.65	0.71	0.72	0.74

## 7. Related Work

To the best of our knowledge, the closest work to the correlation miner has been published in <sup>12</sup>, where the authors tackled the same problem, but based their analysis on event logs neither with case id nor with timestamps. Since the correlation miner also uses timestamps, it is expected to produce more reliable results. Specifically, for the real-world event log (Sec. 4.3), the algorithm presented in <sup>12</sup> yields a model with precision 61% and recall 41%, which are 24% and 22% lower than the correlation miner's results respectively. However, <sup>12</sup> can process cyclic orchestration models.

Correlating events in an event log is a continuing challenge in the field of process mining, and more specifically, as suggested in <sup>19</sup>, this challenge has a high importance in service mining. Therefore, we discuss related work in three categories: (i) process discovery, (ii) service mining, and (iii) event correlation.

A large and growing body of literature has been published in the field of process discovery. Van Dongen et al. have reviewed discovery algorithms that generate process models in the Petri Net notation <sup>2</sup>. Also, other surveys such as the one conducted by Tiwari et al. <sup>7</sup>, have shown that a significant number of studies employ different approaches to discover process models from event logs, including, but not limited to, data mining techniques <sup>20,21</sup>, genetic algorithms <sup>22,23</sup>, fuzzy algorithms <sup>3</sup>, markovian approaches <sup>24</sup> and cluster analysis <sup>25</sup>. A recent study by Verbeek and van der Aalst <sup>8</sup>, in which the authors have introduced a generic divide-and-conquer framework to enable process discovery and conformance checking for big event logs. Their approach is similar to our correlation miner in the sense that both have been implemented using Integer Linear Programming. However, the precise optimization problem that they solve is different.

The topic of service mining has been investigated in <sup>6</sup>, where the authors have illustrated the potential of applying process mining in the context of web services by employing the ProM framework as a process mining tool and IBM's WebSphere as a reference system. Furthermore, Dustdar et al. <sup>26</sup> have introduced Web Services Interaction Mining, which concerns itself with performing process mining techniques in order to analyze service interactions. In <sup>27</sup> the authors have presented a web service mining framework aiming at the discovery of unexpected and interesting service compositions that automatically emerge in a bottom-up fashion. Finally, a study by van der Aalst <sup>19</sup> has reported challenges in service mining, which include the correlation challenge that we have tackled in this paper.

To the best of our knowledge, apart from <sup>12</sup>, no process or service mining algorithms currently exist that can mine a log that contains no case identifiers and no additional data elements to do the correlation on.

Several techniques have been developed to facilitate event correlation in the context of service-oriented systems. <sup>28</sup> identified a set of 18 correlation patterns that have been grounded in a formal model. The concept of correlation set, a query to retrieve identifiers from messages that are unique for a particular cases, has been included in many correlation techniques such as <sup>29</sup>, in which an algorithm has been proposed to assign a certain identifier to each case in a multi-party supply chains. De Pauw et al. carried out a study to discover conversations in web services by using semantic correlation analysis <sup>30</sup>, in which different services pass dedicated identifiers inside their messages. Moreover, in <sup>31</sup> the authors developed an interactive semi-automated tool for event correlation from web service interaction logs.

In contrast, our algorithm enables fully automated event correlation and is only based on occurrences and timestamps of events rather than other data elements that are associated with events. In particular, our technique can be used to mine service collaboration scenarios, in which no main controller exists to manage interactions among the services, and therefore, no data element such as payload of services can be recorded (e.g., <sup>32</sup>) whereas most of the related studies in which correlation will be performed on some data elements that has been recorded by a main controller.

## 8. Conclusion

In this paper, we have presented a process discovery technique, the correlation miner, which enables mining of both business process models and causal relations from event logs that do not contain case identifiers. The correlation miner consists of two phases.

The core idea for the first phase of the correlation miner is to construct a business process model, in such a way that the number of cases that flow into and out of an activity are equal to the number of events that happen for that activity. However, it is often possible to generate more than one process model that meets this rule. Therefore, we defined additional criteria that need to be fulfilled in order to select the best model. Since logs without case identifiers provide two elements of information: (i) how many times an event occurred for a particular activity, and (ii) at which time an event has occurred, we designed our algorithm based on these two characteristics.

Accordingly, the process discovery part of correlation miner has three steps. In the first two steps, it creates two matrices, the Precede/Succeed matrix and the Duration matrix. Given two activities, the corresponding element in the Precede/Succeed matrix indicates the fraction of events referring to the first activity that have occurred before events referring to the second activity. If this value is high, it is more likely that there is an edge from the first to the second activity. Similarly, given two activities, the corresponding element in the Duration matrix indicates the average time difference between events referring to the first activity and events referring to the second activity. If this value is low, it is more likely that there is an edge from the first to the second activity. Finally, in the third step, the correlation miner constructs all possible business process models that meet the rule mentioned above and then it selects the best one based on the values from the Precede/Succeed matrix and the Duration matrix.

The second phase of the correlation miner aims at reconstructing cases from the event log. The algorithm proposed for this phase, works by interpreting an event log as a set of binary relations between events, where a relation from one event to another means that this event is directly followed by another event in the same case. To determine these relations, the correlation miner tries to cast the problem into a Quadratic Programming (QP) problem, in which a set of constraints must hold for the binary relations, and an optimization function selects the 'best' set of binary relations. This objective function is designed to comply with the assumption that time differences between events that belong to the same case are more likely to follow a certain probability distributions, than time differences between binary relations from different cases.

To evaluate the applicability of the correlation miner, we performed experiments with both synthetic event logs and a real-world event log. The results from the evaluation show that the correlation miner produces good results under most conditions. In particular, it produced a business process model with 85% precision

and 63% recall for the real-world log with 11,647 cases. A set of truncated versions of this real-world event log was used further to evaluate the second phase of correlation miner, which resulted in a reconstructed event log with 74% precision and 71% recall when outliers were removed from the log. Furthermore, the evaluations with the synthetic logs show that results are better for structured models than for unstructured models. Also, results are better – specially for the first phase – when there are more cases in the log to mine from, and when the inter-arrival time of cases is higher than the duration of activities.

The current correlation miner is only able to mine acyclic business processes. There is, therefore, a definite further need for research to enable mining of cyclic processes as well. Furthermore, additional work is needed to ensure that our algorithm can produce an accurate average time difference between activities for the duration matrix, given that we do not know which events belong to the same case and should, therefore, be used to compute the time difference. As a closing note, there is room for further improvement in determining more efficient mechanism to explore the constraints in the second phase of correlation miner, which may improve the time performance of this algorithm when event logs with very large numbers of cases are mined.

### Acknowledgments

The research leading to these results has received funding from the European Union's Seventh Framework Programme under grant agreement 2012-318275 (GET Service).

### References

1. W. Van Der Aalst, A. Adriansyah, A. K. A. de Medeiros, F. Arcieri, T. Baier, T. Blickle, J. C. Bose, P. van den Brand, R. Brandtjen, J. Buijs, et al. Process mining manifesto. In *Business process management workshops*, pages 169–194. Springer, 2011.
2. B. F. Van Dongen, A. A. De Medeiros, and L. Wen. Process mining: Overview and outlook of petri net discovery algorithms. In *Trans. on Petri Nets and Other Models of Concurrency II*, pages 225–242. Springer, 2009.
3. C. W. Günther and W. M. van der Aalst. Fuzzy mining–adaptive process simplification based on multi-perspective metrics. In *BPM*, pages 328–343. Springer, 2007.
4. A. Rozinat and W. M. van der Aalst. Conformance checking of processes based on monitoring real behavior. *Information Systems*, 33(1):64–95, 2008.
5. W. van der Aalst. *Process Mining*. Springer, Heidelberg - Berlin, 2011.
6. W. M. van der Aalst and H. E. Verbeek. Process mining in web services: The websphere case. *IEEE Data Eng. Bull.*, 31(3):45–48, 2008.
7. A. Tiwari, C. J. Turner, and B. Majeed. A review of business process mining: state-of-the-art and future trends. *BPM Journal*, (1):5–22, 2008.
8. H. Verbeek and W. van der Aalst. Decomposed process mining: The ilp case. In *BPI 2014 Workshop, Haifa, Israel (September 2014) Accepted for publication*, 2014.
9. C. Günther and A. Rozinat. Disco: Discover your processes. In *Proceedings of the BPM 2012 Demo Track*, volume 940 of *CEUR Workshop Proceedings*, pages 40–44, 2012.

10. W. van der Aalst. Service mining: Using process mining to discover, check, and improve service behavior. *IEEE Trans. on Services Computing*, 6(4):525–535, 2013.
11. S. Pourmirza, R. Dijkman, and P. Grefen. Correlation mining: Mining process orchestrations without case identifiers. In *Service-Oriented Computing*, pages 237–252. Springer, 2015.
12. D. R. Ferreira and D. Gillblad. Discovering process models from unlabelled event logs. In *Business Process Management*, pages 143–158. Springer, 2009.
13. D. B. Johnson. Finding all the elementary circuits of a directed graph. *SIAM Journal on Computing*, 4(1):77–84, 1975.
14. G. Optimization. Inc.. gurobi optimizer reference manual, version 5.0, 2012.
15. B. F. van Dongen, A. K. A. de Medeiros, H. Verbeek, A. Weijters, and W. M. Van Der Aalst. The prom framework: A new era in process mining tool support. In *Applications and Theory of Petri Nets 2005*, pages 444–454. Springer, 2005.
16. J. Mendling, H. A. Reijers, and J. Cardoso. What makes process models understandable? In *BPM*, pages 48–63. Springer, 2007.
17. B. van Dongen. Bpi challenge 2012. dataset (2012).
18. B. Efron. Nonparametric estimates of standard error: the jackknife, the bootstrap and other methods. *Biometrika*, 68(3):589–599, 1981.
19. W. M. van der Aalst. Challenges in service mining: record, check, discover. In *Web Engineering*, pages 1–4. Springer, 2013.
20. M. Hammori, J. Herbst, and N. Kleiner. *Interactive workflow mining*. Springer, 2004.
21. S.-Y. Hwang and W.-S. Yang. On the discovery of process models from their instances. *Decision Support Systems*, 34(1):41–57, 2002.
22. W. M. Van der Aalst, A. A. de Medeiros, and A. Weijters. Genetic process mining. In *Applications and theory of Petri nets 2005*, pages 48–69. Springer, 2005.
23. A. de Medeiros and A. Karla. Genetic process mining, 2006.
24. H. Mannila and D. Rusakov. Decomposition of event sequences into independent components. In *First SIAM Int. Conf. on Data Mining*. SIAM, 2001.
25. G. Greco, A. Guzzo, L. Pontieri, and D. Sacca. Mining expressive process models by clustering workflow traces. In *Advances in Knowledge Discovery and Data Mining*, pages 52–62. Springer, 2004.
26. S. Dustdar and R. Gombotz. Discovering web service workflows using web services interaction mining. *IJBPM*, 1(4):256–266, 2006.
27. G. Zheng and A. Bouguettaya. Service mining on the web. *Services Computing, IEEE Trans. on*, 2(1):65–78, 2009.
28. A. Barros, G. Decker, M. Dumas, and F. Weber. Correlation patterns in SoA. In *Fundamental Approaches to Software Engineering*, pages 245–259. Springer, 2007.
29. K. Gerke, J. Mendling, and K. Tarmyshov. Case construction for mining supply chain processes. In *Business Information Systems*, pages 181–192. Springer, 2009.
30. W. De Pauw, R. Hoch, and Y. Huang. Discovering conversations in web services using semantic correlation analysis. In *ICWS*, pages 639–646. IEEE, 2007.
31. H. R. Motahari-Nezhad, R. Saint-Paul, F. Casati, and B. Benatallah. Event correlation for process discovery from web service interaction logs. *The VLDB Journal*, 20(3):417–444, 2011.
32. S. Pourmirza, R. Dijkman, and P. Grefen. Switching parties in a collaboration at run-time. In *Enterprise Distributed Object Computing Conference (EDOC), 2014 IEEE 18th International*, pages 136–141. IEEE, 2014.