

Mining Context-Dependent and Interactive Business Process Maps using Execution Patterns

Jiafei Li^{1,2}, R.P. Jagadeesh Chandra Bose², and Wil M.P. van der Aalst²

¹ College of Computer Science and Technology, Jilin University, China 130012,
jiafei@jlu.edu.cn, l.j.f.li@tue.nl,

² Eindhoven University of Technology, The Netherlands
j.c.b.rantham.prabhakara@tue.nl, w.m.p.v.d.aalst@tue.nl

Abstract. Process mining techniques attempt to extract non-trivial knowledge and interesting insights from event logs. Process models can be seen as the “maps” describing the operational processes of organizations. Unfortunately, traditional process discovery algorithms have problems dealing with less-structured processes. Furthermore, existing discovery algorithms do not consider the analyst’s context of analysis. As a result, the current models (i.e., “maps”) are difficult to comprehend or even misleading. To address this problem, we propose a two-phase approach based on common execution patterns. First, the user selects relevant and context-dependent patterns. These patterns are used to obtain an event log at a higher abstraction level. Subsequently, the transformed log is used to create a hierarchical process map. The approach has been implemented in the context of ProM. Using a real-life log of a housing agency we demonstrate that we can use this approach to create maps that (i) *depict desired traits*, (ii) *eliminate irrelevant details*, (iii) *reduce complexity*, and (iv) *improve comprehensibility*.

1 Introduction

Process mining aims at extracting process-related information from event logs. Process mining techniques can deliver valuable, factual insights into how processes are being executed in real life. The majority of research in process mining so far has focussed on process discovery (both from a control-flow and organizational perspective). Process models can be seen as the “maps” describing the operational processes of organizations. Unfortunately, *accurate and interactive business process maps* are missing. Either there are no good maps or maps (if available) are static and/or outdated [1].

Process mining techniques can be used to generate process maps [2–4]. We have applied our process mining tool ProM in more than 100 organizations and our experiences show that processes tend to be less structured than expected. Traditional process discovery algorithms have problems dealing with such unstructured processes and generate spaghetti-like process models that are hard to comprehend. The granularity at which the events are logged is typically different from the desired granularity. Analysts and end users prefer a higher level of abstraction without being confronted with lower level events stored in raw event logs.

Analogous to cartography, process mining techniques should allow for various context-dependent views on the process maps. For example, the perspective of analysis may be different depending on someone's role and expertise e.g., a manager may be interested in a high level view, while a specialist may be interested in a detailed analysis of some process fragment. Process discovery techniques should facilitate the extraction of process maps eliciting the respective desired traits and hiding the irrelevant ones for various users. Furthermore, these techniques should uncover comprehensible models by providing a hierarchical view with a facility to seamlessly zoom in or zoom out the process maps. There is an imperative need for *techniques that automatically generate understandable and context-dependent business process maps* [1].

In this paper, we propose a *two-phase approach to mine interactive and context-dependent business process maps based on common execution patterns*. The *first phase* comprises the pre-processing of log with desired traits and at a desired level of granularity. This paper will show one means to realize this by uncovering common execution patterns in the log, selecting context-dependent patterns, and defining abstractions over these patterns. Pattern selection and the mapping with abstractions can be interactively performed by the user. Event logs are then pre-processed (transformed) with these abstractions. In the *second phase*, the transformed log is used for process discovery. Any discovery algorithm with an ability to zoom-in/out the sub-processes defined by the abstractions can be used. This paper presents an adapted version of the Fuzzy Miner [3] and shows that it can provide such hierarchical view of process maps. The two-phase approach presented in this paper has been implemented in ProM 6.0³. Figure 1 highlights the difference between the traditional approach to do process discovery and the two-phase approach. Note that the process model (map) mined using the two-phase approach is simpler and that this approach enables the abstraction of activities based on functionality and provides a seamless zooming into the sub-processes captured in the abstractions.

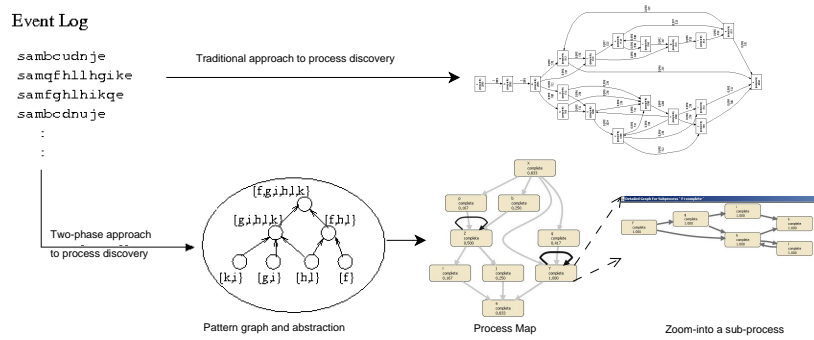


Fig. 1: Traditional approach vs. two-phase approach

³ ProM 6.0 is not officially released yet, but nightly builds, including the reported functionality are available from www.processmining.org

The remainder of the paper is organized as follows. Our two-phase approach to mining process maps is introduced in Section 2. Section 3 presented pattern definitions and pattern metrics while Section 4 proposes one approach to form abstractions based on patterns. In Section 5, we detail our two-step approach and describe an adaptation of Fuzzy Miner to discover process maps. Section 6 presents a case study of a real-life log from a rental agency. Related work is discussed in Section 7. Section 8 concludes the paper.

2 Two-phase Approach to Mine Process Maps

We use the following notations in this paper. Let Σ denote the set of activities. $|\Sigma|$ is the number of activities. Σ^+ is the set of all non-empty finite sequences of activities from Σ . We denote traces by bold face lower case letters \mathbf{t}_1 , \mathbf{t}_2 etc. A trace \mathbf{t} is an element of Σ^+ . $\mathbf{t}(i)$ denotes the i^{th} activity in the trace. For $i < j$, $\mathbf{t}(i, j)$ denotes the subsequence from the i^{th} position to the j^{th} position in the trace \mathbf{t} . An event log \mathcal{L} corresponds to a bag (i.e., a multiset) of traces.

Phase-1: Preprocessing Log In this phase, the log is simplified based on the desired traits of the context of analysis. A mapping $\mathcal{M} \subseteq 2^\Sigma \times \mathcal{A}$ is defined between the *original alphabet* of the event log Σ , and an *abstract alphabet* \mathcal{A} . An example mapping is $\mathcal{M} = \{(\{\mathbf{a}, \mathbf{b}\}, \mathbf{x}), (\{\mathbf{b}, \mathbf{c}, \mathbf{d}\}, \mathbf{y}), (\{\mathbf{e}\}, \mathbf{z}), (\{\mathbf{d}\}, \mathbf{z})\}$. This mapping is analogous to the grouping and tagging of streets as a town/city in cartography and to the selection of a desired perspective of viewing maps (restaurant maps vs. fuel station maps). The analyst can define this mapping based on domain knowledge or can be assisted by uncovering common execution patterns and relationships between them in the log. These common execution patterns typically capture a sub-process/functionality. Analysts would like to capture such subprocess behavior in its totality as an *abstract activity* in a mined process model with a facility to zoom in/out the subprocess if needed. The mapping is defined over the sets of activities manifested as patterns. We present techniques that assist in automatically uncovering such patterns and relationships between activities in Section 3.

$\mathcal{D} = \bigcup_{(A, \mathbf{a}) \in \mathcal{M}} A$ denotes the set of activities in Σ for which a mapping is defined. The original event log \mathcal{L} , is transformed into an *abstract log* \mathcal{L}' . Each trace $\mathbf{t} \in \mathcal{L}$ is transformed into a corresponding trace $\mathbf{t}' \in \mathcal{L}'$. In each trace \mathbf{t} , the manifestation of each pattern captured by $(A, \mathbf{a}) \in \mathcal{M}$ is replaced with its *abstract activity*, \mathbf{a} , in the transformed trace. The activities in $\Sigma \setminus \mathcal{D}$ being not involved in the definition of mapping indicate activities that are insignificant from the context of analysis and are filtered from \mathbf{t} during this transformation. In Section 5, we describe the transformation of log in detail.

Phase-2: Mining Maps The second phase is to mine a process model on the abstract log. The mapping defined in Phase-1 induces a hierarchy over the abstract activities. Upon zooming into an abstract activity, a process model depicting the subprocess captured by this abstract activity is shown. The patterns replaced by the abstract activity are used to create this sub-process model. We

adapted Fuzzy Miner for this phase and the details are presented in Section 5. Note that this is a generic approach that can be iterated over any number of times with the event log for iteration $i + 1$ being the output event log of iteration i .

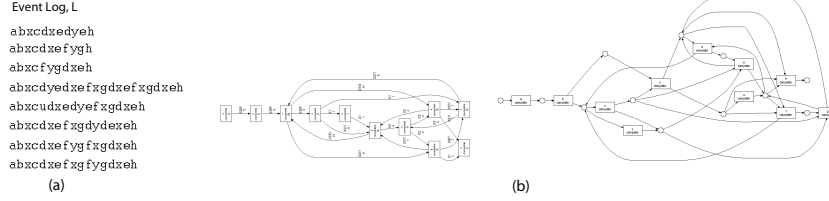


Fig. 2: (a) An example log (b) process models discovered using Heuristic miner and Alpha++ miner

We use a running example log depicted in Figure 2(a) to illustrate the approach. This log contains 8 process instances with 11 event classes. Figure 2(b) depicts two models mined using traditional process discovery techniques. It is imperative to find that both these models are not easy to understand. In the following sections, we will present our two-phase approach in more detail using this example.

3 Pattern Definitions and Pattern Metrics

In this section, we adapt the pattern definitions proposed in [5] and define metrics over these patterns. We consider only the *maximal repeat* patterns for the discussion in this paper. However, other patterns such as *tandem arrays* capturing the manifestation of loop constructs proposed in [5] can also be used. These patterns are later used to define the mapping \mathcal{M} between activities and abstractions.

3.1 Pattern definitions

Definition 1 (Maximal Repeat). A maximal pair in a sequence, s , is a pair of identical sub-sequences α and β such that the symbol to the immediate left (right) of α is different from the symbol to the immediate left (right) of β . In other words, extending α and β on either side would destroy the equality of the two strings. A maximal pair is denoted by the triple (i, j, α) where i and j correspond to the starting positions of α and β in s with $i \neq j$. A maximal repeat in a sequence, s , is defined as a subsequence α that occurs in a maximal pair in s .

Maximal repeats capture execution patterns (sequence of activities) common within a trace and/or across a set of traces in an event log. Such patterns might be evidence of common functionality (often abstracted as a sub-process). In order to find these commonalities across multiple traces in the entire event log, we first construct a single sequence, say, s , which is obtained by the concatenation of traces in the event log with a distinct delimiter between the traces. Maximal repeats are then discovered over this concatenated sequence s . Maximal repeats

can be efficiently discovered in linear time using suffix trees for strings [6]. Let $\mathcal{P}_{\mathcal{L}}$ be the set of all patterns in $\log \mathcal{L}$. In this paper, $\mathcal{P}_{\mathcal{L}}$ includes all maximal repeats and $\{a \mid a \in \Sigma\}$. A *base pattern* is a pattern that does not contain any other pattern within it. The pattern **abxc** is a base pattern while **abxcdxe** is not because the latter pattern contains the pattern **x** within it. Let $\mathcal{P}_{\mathcal{L}}^b$ be the set of all base patterns in $\mathcal{P}_{\mathcal{L}}$.

Consider the trace $\mathbf{t}_6 = \mathbf{abxcdxefxgdyedxe}$ in the log of Figure 2(a). The maximal pairs in \mathbf{t}_6 are (3, 6, **x**), (5, 11, **d**), (7, 13, **e**) and (5, 14, **dx**). There are a total of 39 maximal repeats in the example log (e.g., **abxc**, **abxcd**, **abxcdxe**, **dx**, **fxg**, **dye**, **fyg**, **dyedxe**, **dxedye**, **h**). $\mathcal{P}_{\mathcal{L}}^b = \{a, b, c, u, x, d, g, f, h, e, y, gf, gd, ef, ed, eh, fxg, dye, fyg, dx, fxgd, abxc, gdxe, dxeh, dxef, efxg, efxgd, abxcd, gdxeh, fygdxe\}$.

Definition 2 (Pattern Alphabet). *The pattern alphabet $\Gamma(p)$, of a pattern, $p \in \mathcal{P}_{\mathcal{L}}$, is the set of activities that appear in p .*

Definition 3. *The equivalence class of a pattern alphabet PA , is defined as $[PA] = \{p \mid p \text{ is a pattern and } \Gamma(p) = PA\}$*

For example, for the patterns **fxg**, **dyedxe** and **dxedye**, the pattern alphabets correspond to $\{f, x, g\}$, $\{d, x, y, e\}$, and $\{d, x, y, e\}$ respectively. The equivalence class of the pattern alphabet $\{d, x, y, e\}$ is $\{dyedxe, dxedye\}$. Equivalence classes of pattern alphabets capture variations of patterns e.g., due to parallelism.

3.2 Pattern Metrics

Pattern metrics such as the frequency of occurrence, significance etc. need to be estimated. A careful consideration needs to be done when estimating the frequency of a pattern. It is due to the fact that certain regions in a trace can contribute to more than one pattern (in the case of overlapping patterns) and might result in misleading frequency counts. For example, consider the trace $\mathbf{t} = \mathbf{abxcdxedfxgdxeh}$ and the pattern alphabet equivalence classes $[\{a, b, x, c\}] = \{\mathbf{abxc}\}$, $[\{a, b, x, c, d\}] = \{\mathbf{abxcd}\}$, and $[\{d, x, e\}] = \{\mathbf{dx}, \mathbf{dxed}\}$. Now, *what should be the pattern (alphabet) counts?* If we consider each of the patterns separately, the pattern, frequency-count pairs are $(\mathbf{abxc}, 1)$, $(\mathbf{abxcd}, 1)$, $(\mathbf{dx}, 2)$ and $(\mathbf{dxed}, 1)$. If we define the pattern alphabet count to be the sum of counts of the patterns captured in its equivalence class, then the pattern alphabet, frequency-count pairs are $(\{a, b, x, c\}, 1)$, $(\{a, b, x, c, d\}, 1)$, and $(\{d, x, e\}, 3)$. It is imperative to see that certain regions in the trace are contributing to more than one pattern (alphabet). The activities in the subsequence $\mathbf{t}(5, 7)$ contributed to two patterns viz., **dx** and **dxed**. Similarly, the activities in the subsequence $\mathbf{t}(1, 4)$ contributed to two patterns viz., **abxc** and **abxcd**.

We identify three distinct methods of dealing with overlaps and counting pattern occurrences. The above method of computing pattern (alphabet) counts is referred to as *Overlapping Alphabet Count (OAC)*. The significance computed using overlapping alphabet counts may be misleading. A more accurate method of computing pattern frequencies is to consider non-overlapping pattern counts. We

distinguish two variations here: (i) considering non-overlap counts for each alphabet separately (local) and (ii) considering non-overlap counts across all alphabets (global) in the event log. These two metrics are referred to as *Non-Overlapping Alphabet Count (NOAC)* and *Non-Overlapping Global Alphabet Count (NOGAC)* respectively. Conflicts arise when more than one pattern can potentially contribute to the count at a region in a trace. One can assign preference to say shorter (longer) patterns to resolve such conflicts. The *NOAC* (with preference to shorter patterns) for the above example is $(\{a, b, x, c\}, 1)$, $(\{a, b, x, c, d\}, 1)$, and $(\{d, x, e\}, 2)$. Note that the conflict at position 5 in **t** for pattern alphabet $\{d, x, e\}$ is resolved in favor of the pattern **dx_e** thereby making **t(5,7)** contribute to only one pattern. The *NOGAC* across all alphabets (with preference to longer patterns) is $(\{a, b, x, c\}, 0)$, $(\{a, b, x, c, d\}, 1)$, and $(\{d, x, e\}, 1)$. A position/subsequence in a trace can contribute to more than one pattern alphabet when considering *NOAC* for each alphabet separately (e.g., index 1 in **t**) while in *NOGAC*, a position contributes to at most one pattern alphabet.

In order to assess the significance of a pattern alphabet PA , we define a metric *Conservedness* (CON_{PA}) = $\frac{NOAC}{\mu} * (1 - \frac{\sigma}{\mu}) * 100\%$ where μ and σ are the mean and standard deviation of the frequencies of activities in PA . *Conservedness* measures the degree to which the individual activities involved in the pattern alphabet manifest as the patterns defined by the alphabet. For example, consider the non-overlap alphabet count of three pattern alphabets $(\{d, x, e\}, 100)$, $(\{d, x, e, f\}, 60)$, and $(\{d, x, e, h\}, 40)$. Let the frequency of activities be $(d, 100)$, $(x, 100)$, $(e, 100)$, $(f, 60)$, and $(h, 40)$. Conservedness value of the pattern alphabets $\{d, x, e\}$, $\{d, x, e, f\}$, and $\{d, x, e, h\}$ is 100%, 51% and 30% respectively. The formal definitions of the above pattern metrics are presented in Appendix A.

4 Abstractions based on Patterns

4.1 Pattern Graph

Relationships exist between patterns (alphabets). For example, consider the patterns **dxefxg**, **dx_e**, and **fxg**. It could be the case that **dx_e** and **fxg** are sub-functionalities used also in a larger context **dxefxg**. One can try to define a partial order capturing the relationships on the pattern alphabets. For example, subsumption can be used as the cover relation. A pattern alphabet PA_i is defined to cover another pattern alphabet PA_j if $PA_j \subset PA_i$ and there is no PA_k such that $PA_j \subset PA_k \subset PA_i$. A *pattern graph* $G = (V, E)$, is a Hasse diagram defined over the partial order on the pattern alphabets, where $V = \{PA_1, PA_2, \dots, PA_n\}$ represents the set of pattern alphabets and E denotes the set of edges (PA_i, PA_j) defined by the cover relation. One can choose either $\mathcal{P}_{\mathcal{L}}$ or $\mathcal{P}_{\mathcal{L}}^b$ to define V . Figure 3(a) depicts a pattern graph on some of the pattern alphabets identified for the example log. We considered pattern alphabets defined by $\mathcal{P}_{\mathcal{L}}^b$ with a conservedness value above 17% to generate this graph.

4.2 Pattern Selection

Nodes in a pattern graph form the basis for abstraction. An analyst can select the pattern nodes based on domain knowledge or by using the pattern metrics

defined in Section 3. We provide two types of node selection modes for abstraction.

Single Node Mode: All manifestations of patterns under the equivalence class of this node’s pattern alphabet are represented by the same abstract activity in the transformed log.

Sub-graph Mode: All manifestations of patterns under the equivalence classes of the pattern alphabets defined by the *induced subgraph* at the selected node are substituted by the abstract activity of the selected node during transformation.

It could be the case that a pattern graph contains a large number of nodes. We recommend to first filter the nodes in the pattern graph before considering them for abstractions. All the metrics defined in Section 3.2 can be used to prune the graph. For example, consider the pattern alphabets $\{a,b,x,c\}$ and $\{a,b,x,c,d\}$ in Figure 3(a). The *NOGAC* of $\{a,b,x,c,d\}$ with preference to shorter patterns (ignoring individual activity patterns) is zero. Similarly, the *NOGAC* of $\{d,x,e,f\}$, $\{d,x,e,h\}$, $\{g,d,x,e,h\}$, $\{g,d,x,e\}$, $\{e,h\}$, $\{e,d\}$, $\{e,f\}$, $\{g,d\}$ and $\{g,f\}$ are all zero. This indicates that manifestations of all patterns under the equivalence class of these pattern alphabets in the log are overlapping with some other pattern. For example, the equivalence class of the pattern alphabet $\{e,d\}$ is $\{ed\}$. There are two manifestations of the pattern ed in \mathcal{L} (in traces $abxc dxedyeh$ and $abxcudxyefxgd xeh$). However, both of these manifestations overlap with $dx e$ and $dy e$ in the example log; thus making the *NOGAC* of $\{e,d\}$ as 0.

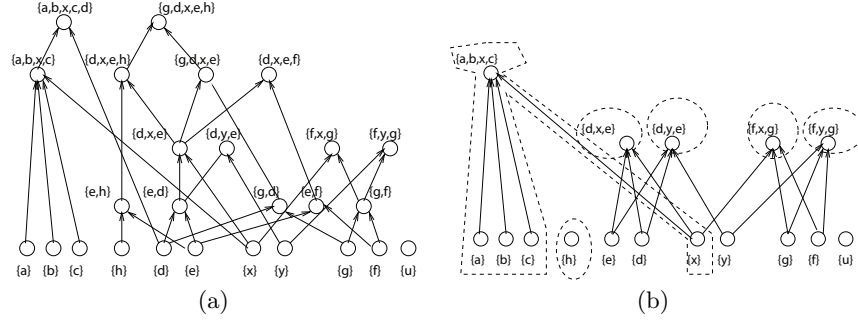


Fig. 3: (a) Pattern graph (b) pattern graph with abstractions for the example log

We recommend to consider nodes capturing longer patterns with a high conservedness value and significant *NOAC* and *NOGAC* to be used under sub-graph mode for abstractions. However, for two pattern alphabet nodes PA_i and PA_j such that $(PA_i, PA_j) \in E$ (i.e., $PA_i \subset PA_j$), if $CON_{PA_i} > CON_{PA_j}$ then, we recommend to consider PA_i under sub-graph mode instead of PA_j though PA_i captures shorter patterns. For example, consider the pattern alphabets $\{d,x,e\}$, $\{d,x,e,f\}$, $\{d,x,e,h\}$ and $\{g,d,x,e\}$. The conservedness value for these alphabets are 52%, 22%, 21%, 22% respectively. It could be seen that the pattern $dx e$ (defined by the alphabet $\{d,x,e\}$) occurs in different contexts in \mathcal{L} . The different contexts are captured by the other three alphabets and are reflected with the relatively low conservedness values for these three alphabets. We recommend to consider $\{d,x,e\}$ as a node for abstraction instead of the other three. Coinciden-

tally in this example the *NOGAC* (with preference to shorter patterns) for the three larger alphabets is also zero.

If nodes in the sub-graph of a pattern node PA_i are covered by one or more nodes PA_j not in the sub-graph of PA_i then we recommend to consider PA_i under single-node mode for abstraction (assuming PA_i is selected). For example, the node $\{d, x, e\}$ is recommended to be considered under single-node mode because the nodes $\{d\}$ and $\{e\}$ in the sub-graph of $\{d, x, e\}$ is also covered by another node $\{d, y, e\}$. Note that these are just recommendations and an analyst can make exceptions if it makes sense according to the context of analysis. Using these guidelines we use the abstractions as defined in Figure 3(b). Here $\{a, b, x, c\}$ is used in the sub-graph mode while $\{d, x, e\}$, $\{d, y, e\}$, $\{f, x, g\}$, $\{f, y, g\}$ and $\{h\}$ are chosen under single-node mode. Certain nodes not pertaining to the context of analysis can also be filtered out (e.g., $\{u\}$). Let us define the mapping \mathcal{M} as $\{(\{a, b, x, c\}, A1), (\{a\}, A1), (\{b\}, A1), (\{c\}, A1), (\{x\}, A1), (\{d, x, e\}, A2), (\{f, x, g\}, A3), (\{d, y, e\}, A4), (\{f, y, g\}, A5), (\{h\}, A6)\}$ on the abstractions chosen for the example log.

5 Process Discovery based on Patterns

5.1 Transformation of Log

Algorithm 1 presents the detail of transforming the log based on the patterns. The *basic idea is to first replace the continuous and intermittent manifestation of each pattern alphabet chosen for abstraction with its abstract activity and make the corresponding low level manifestations part of the sub-log corresponding to the abstract activity. The sub-log of an abstract activity can be used to zoom in the detailed behavior.* The intermittent manifestation here refers to the situation where the execution of the subsequence corresponding to a pattern is interrupted by other activities. For instance, let *dye* be a pattern, the manifestation of *dye* in the trace *abxcdxrefxgdydexeh* is intermittent because *dye* is interrupted by *d*.

```

1: Let  $\mathcal{M}$  be the mapping chosen by the user.  $\mathcal{A} = \cup_{(PA, a) \in \mathcal{M}} \{a\}$  defines the set of defined
   abstractions.  $\mathcal{SP} = \cup_{(PA, a) \in \mathcal{M}} [PA]$  denotes the set of all patterns for which
   abstractions are defined. Let  $f: \mathcal{SP} \rightarrow \mathcal{A}$  be the function defining the abstraction for
   each pattern. Let  $l: \mathcal{A} \rightarrow \mathcal{SL}$  be the function defining the sub-log for each abstraction.
   Let  $\mathcal{L}'$  be the transformed log of  $\mathcal{L}$ . Initialize  $\mathcal{L}' = \{\}$  and  $l(a) = \{\}$  for all  $a \in \mathcal{A}$ 
2: for all  $t \in \mathcal{L}$  do
3:   Let  $t'$  be an empty trace. Set  $j = 1$ .
4:   while  $j \leq |t|$  do
5:     Let  $LD_s$  be the list of patterns in  $\mathcal{SP}$  starting with  $t(j)$  ordered in descending
       order of their length
6:     for every pattern  $\alpha \in LD_s$  do
7:       if there exists a continuous manifestation of a pattern  $\alpha$  at index  $j$  in  $t$  then
8:          $l(f(\alpha)) = l(f(\alpha)) \cup \{t(j, j + |\alpha|)\}$ ; Append  $f(\alpha)$  to  $t'$ ; Set  $j = j + |\alpha| - 1$ ; exit for
9:       else if there exists an intermittent manifestation of  $\alpha$  at index  $j$  in  $t$  then
10:        Re-adjust the intermittent manifestation in  $t$ .
11:         $l(f(\alpha)) = l(f(\alpha)) \cup \{\alpha\}$ ; Append  $f(\alpha)$  to  $t'$ ; Set  $j = j + |\alpha| - 1$ ; exit for
12:      end if
13:    end for
14:    Set  $j = j + 1$ 
15:  end while
16:   $\mathcal{L}' = \mathcal{L}' \cup \{t'\}$ 
17: end for

```

Algorithm 1: Single-phase pattern-based log transformation

Steps 9-12 in the algorithm deal with the intermittent manifestation of a pattern and substitutes it with the abstract activity. Algorithm 1 will transform the trace `abxcdxexgdydexeh` in our example log to `A1A2A3A4A2A6`. In this way, one can cope with situations where a common functionality is interrupted by other activities in concurrency.

5.2 Adapting Fuzzy Miner to Discover Maps

ProM's *Fuzzy Miner* [3] is inspired by cartography to provide business process maps. However, the existing miner has some limitations. It (i) cannot customize maps from a defined context (city maps vs. highway maps) (ii) introduces the risk of aggregating unrelated activities together in a cluster (a street in Eindhoven is clustered along with streets in Amsterdam) and (iii) provides two level hierarchy instead of a multi-level hierarchical view of the process map.

We adapted Fuzzy Miner to support the discovery of process maps. The pattern selection techniques presented in Section 4 facilitate customization from an user's context and getting meaningful abstract activities. By using the sub-log of each abstract activity, we implemented the functionality of zooming in/out the abstract activity and showing the detailed sub-process captured by it. Furthermore, by combining with the existing functions in the Fuzzy Miner of zooming in/out the cluster nodes, a three-level view of the process map is provided.

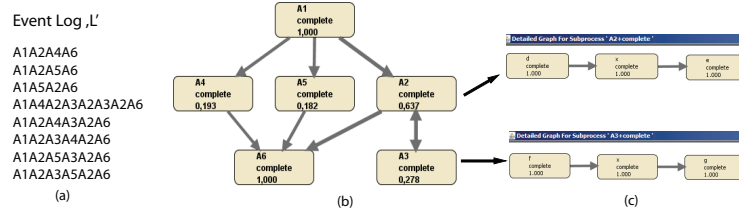


Fig. 4: (a) Transformed Log (b) process map mined from the transformed log using adapted Fuzzy miner (c) sub-process maps when zooming in on A2 and A3

To illustrate this two-phase approach, let us transform the log as described in Algorithm 1 using the mapping \mathcal{M} . The transformed log is shown in Figure 4(a). Figure 4(b) depicts the process map mined by the adapted Fuzzy miner, while Figure 4(c) shows the sub-process maps when zooming in the abstract activities A2 and A3. It is evident from Figure 2 and Figure 4 that our two-phase approach helps in presenting more accurate and more readable models.

6 Case Study and Discussion

We applied the techniques proposed in this paper on a real life log of a rental agency where the cases corresponded to cancellation of a current rental agreement and subsequent registration of a new rental agreement. This log was provided by a large Dutch agency that rents houses and apartments and contains 210 cases, 6100 events and 74 event classes. Figure 5(a) depicts the process model mined using heuristic miner. This process model included two types of

cancellation as highlighted by two rectangles in Figure 5(a). The unselected region corresponds to common functionality used by both of them. The resulting model is difficult to comprehend.

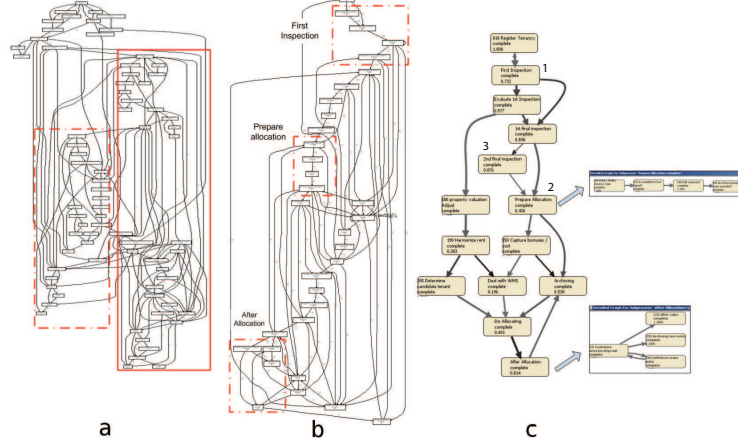


Fig. 5: (a) Heuristic net mined on the whole log (b) heuristic net mined on the filtered log (c) process map and zoomed-in sub-processes mined from transformed log based on interactive and context-dependent abstractions

Table 1: Three of the pattern alphabets chosen for abstraction

No.	Abstraction Name	Pattern Alphabet	<i>NOAC</i>	<i>CON</i> (%)
1	First Inspection	[050 Plans appointment 1st Inspection,060 Edit confirmation letter / Tenancy form, 070 Is 1st inspection performed?,100 Ready report 1st Insp. / Make-Calculation For]	80	66
2	Prepare Allocation	[500 Rate / Modify vacancy type,540 Are there bonuses / costs awarded?, 510 Is completion form signed?530 Edit command]	145	93
3	2nd Final Inspection	[460 (Re) plans 2nd final inspection,470 Modify date lease, 480 is the 2nd final inspection performed?]	4	75

In this study, we concentrate on one type of cancellation process defined by the solid rectangle (let us assume that the analyst wants to focus on this cancellation process). The primary steps involve the registration of a request, multiple inspections of the rented house, determining (future) tenants, (re-)allocation and archiving of the case. We first identified the common execution patterns in this log and chose 17 abstract activities (some involve pattern alphabets and some involve individual activities) concerned with the above primary steps. Table 1 shows some of the pattern alphabets used in defining abstractions. Pattern alphabets capturing a functionality from a domain point of view are chosen as candidate nodes (under sub-graph mode) for abstractions. A meaningful name is defined for every candidate abstraction. Those pattern alphabets with a significant *NOAC* as well as a high *CON* value have priority to be selected for abstractions as can be seen in Table 1. We used these seventeen abstractions to do the first phase of log transformation. Then, to utilize the important pattern of “2nd Final Inspection” which cannot be found in the first iteration, we performed

a second iteration of pattern identification and chose 14 abstract activities for log transformation and pattern discovery. The resulting process map mined by the adapted Fuzzy miner is shown in Figure 5(c) which also presents the sub-process when zooming in the abstract activities of **Prepare Allocation** and **After Allocation**. Each sub-process subsumes the manifestation of patterns captured in the sub-log defined by the abstraction.

To have a fair comparison, we filtered the activities from the original log that were not involved in the abstraction mapping. This filtered log contains 210 instances, 4784 events and 33 event classes. The process model discovered by heuristic miner on the filtered log is shown in Figure 5(b). It is apparent that the process map discovered by our two-step approach is more comprehensible and captures the main steps of this specific type of rental cancellation process. This resulting process map not only facilitates the analyst to get an overview of the whole process, but also makes it easy to seamlessly zoom-in each abstract activity to observe the detailed sub-process. This shows that using our two-step approach indeed leads to better understandable process maps without sacrificing precision.

7 Related Work

Several approaches based on trace clustering [7–9] have been proposed in literature. Trace clustering enables the partitioning of the event log based on coherency of cases. Process models mined from each of the resulting clusters are expected to be simpler than that of the one mined from the entire event log. Greco *et al.* [9] augmented trace clustering with an approach to mine hierarchies of process models that collectively represent the process at different levels of granularity and abstraction. This approach tries to analyze the mined process models (post-processing) for identifying activities that can be abstracted. However, for large complex logs, the mined process models (even after clustering) can be quite spaghetti-like. In contrast, the approach proposed in this paper analyzes the raw traces and defines abstraction (pre-processing) and has the ability to zoom-in hierarchically into the abstract entities. Furthermore, the user has flexibility and control when selecting the abstractions/activities of interest based on his/her context of analysis.

Taking cartography as a metaphor, Günther and Aalst [3] have proposed the fuzzy mining approach to implement process simplification. Less significant activities/edges are either removed or clustered together in the model. However, this approach poses a danger of clustering activities/edges having no domain significance. Polyvyanyy *et al.* [10] have proposed a slider approach for enabling flexible control over various process model abstraction criteria. Approaches such as [10, 3] look at abstraction from the point of retaining highly significant information and discarding less significant ones in the process model where the notion of significance is defined over the (relative-)frequency of occurrence of an entity and not based on the context. In contrast, the approach proposed in this paper looks at abstraction from a functionality/subprocess point of view which performs filtering of activities based on the context of analysis. Our approach

can be used as a preprocessing step for the logs and can be seamlessly integrated with other approaches for abstraction [9, 3] as well as with classical approaches for process discovery such as the heuristic approach in [4].

8 Conclusions and Future Work

This paper presented a two-phase approach to mining business process maps that comprises the pre-processing of log based on desired traits and at a desired level of granularity as a first step and discovering the maps with seamless zoom-in facility as the second step. We discussed one means of realizing this two-phase approach by exploiting the common execution patterns in the event log. Metrics assessing the significance of these patterns and ways of selecting these patterns for abstractions were presented. Our initial results on a few real-life logs show encouraging results. Concurrency in process models adds complexity to the discovery of patterns. As future work, we focus on more real-life applications and improving the robustness of the approach in the context of concurrency.

Acknowledgments This work is supported in part by the NNSF (No.60873149) and the 863 project (No.2006AA10Z245) of China, and in part by EIT, NWO-EW, and STW. R.P.J.C. Bose and W.M.P. van der Aalst are grateful to Philips Healthcare for funding the research in Process Mining.

References

1. van der Aalst, W.M.P.: Challenges in Business Process Mining. *Applied Stochastic Models in Business and Industry* (to appear)
2. van der Aalst, W.M.P.: Using process mining to generate accurate and interactive business process maps. In: *BIS (Workshops)*. Volume 37 of *LNBIP*. (2009) 1–14
3. Günther, C.W., van der Aalst, W.M.P.: Fuzzy Mining - Adaptive Process Simplification Based on Multi-perspective Metrics. In: *Business Process Management (BPM)*. Volume 4714 of *LNCS*. (2007) 328–343
4. Weijters, A., van der Aalst, W.M.P.: Rediscovering workflow models from event-based data using little thumb. *Integrated Computer-Aided Engineering* **10**(2) (2003) 151–162
5. Bose, R.P.J.C., van der Aalst, W.M.P.: Abstractions in Process Mining: A Taxonomy of Patterns. In: *Business Process Management (BPM)*. Volume 5701 of *LNCS*. (2009) 159–175
6. Gusfield, D.: *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*. Cambridge University Press (1997)
7. Bose, R.P.J.C., van der Aalst, W.M.P.: Context Aware Trace Clustering: Towards Improving Process Mining Results. In: *Proceedings of the SIAM International Conference on Data Mining (SDM)*. (2009) 401–412
8. Bose, R.P.J.C., van der Aalst, W.M.P.: Trace Clustering Based on Conserved Patterns: Towards Achieving Better Process Models. In: *Business Process Management Workshops*. Volume 43 of *LNBIP*. (2009) 170–181
9. Greco, G., Guzzo, A., Pontieri, L.: Mining Taxonomies of Process Models. *Data Knowl. Eng.* **67**(1) (2008) 74–102
10. Polyvyanyy, A., Smirnov, S., Weske, M.: Process Model Abstraction: A Slider Approach. In: *Enterprise Distributed Object Computing*. (2008) 325–331

This appendix shows some of the technical details. It will not be included in the final paper.

A Definitions of Pattern Metrics

In this appendix, we give formal definitions for the pattern metrics proposed in Section 3.2.

Definition 4 (Event Log). An event Log \mathcal{L} is a multi-set (bag) of traces (\mathcal{T}, f) , where \mathcal{T} is the set of traces and $f : \mathcal{T} \rightarrow \mathbb{N}_{\geq 1}$ denotes the number of occurrences of a trace $\mathbf{t} \in \mathcal{T}$ in \mathcal{L} .

Definition 5 (Base Pattern). A base pattern is a pattern that does not contain any other pattern within it. In other words, a base pattern is one whose length is the same as that of the size of its alphabet.

The pattern **abxc** is a base pattern while the pattern **abxcdxe** is not. The pattern **abxcdxe** contains the pattern **x** in it. The length of the pattern **abxcdxe** is 7 while the size of its pattern alphabet is 6.

Definition 6 (Pattern indices). Let p be a pattern, and \mathbf{t} be a trace of length l , the indices of the pattern p in trace \mathbf{t} is defined as $\mathcal{I}(p, \mathbf{t}) = \{i \mid 1 \leq i \leq l \wedge p \text{ is a prefix of } \mathbf{t}(i, l)\}$. Let PA be a pattern alphabet. The indices of the pattern alphabet PA in trace \mathbf{t} is defined as $\mathcal{I}(PA, \mathbf{t}) = \bigcup_{p \in [PA]} \mathcal{I}(p, \mathbf{t})$.

Definition 7 (Shortest and longest non-overlapping pattern at index). Let \mathbf{t} be a trace and \mathcal{P} be the set of all patterns. The shortest non-overlapping pattern in \mathcal{P} at index i in trace \mathbf{t} is defined as $\text{shortest}(\mathcal{P}, \mathbf{t}, i) = \{p \in \mathcal{P} \mid i \in \mathcal{I}(p, \mathbf{t}) \wedge \forall p' \in \mathcal{P} (i \in \mathcal{I}(p', \mathbf{t}) \Rightarrow |p| \leq |p'|)\}$. $\text{shortest}(\mathcal{P}, \mathbf{t}, i)$ is either empty or a singleton. Let $|\text{shortest}(\mathcal{P}, \mathbf{t}, i)|$ denote the length of the shortest pattern.

The longest non-overlapping pattern in \mathcal{P} at index i in trace \mathbf{t} is defined as $\text{longest}(\mathcal{P}, \mathbf{t}, i) = \{p \in \mathcal{P} \mid i \in \mathcal{I}(p, \mathbf{t}) \wedge \forall p' \in \mathcal{P} (i \in \mathcal{I}(p', \mathbf{t}) \Rightarrow |p| \geq |p'|)\}$. $\text{longest}(\mathcal{P}, \mathbf{t}, i)$ is either empty or a singleton. Let $|\text{longest}(\mathcal{P}, \mathbf{t}, i)|$ denote the length of the longest pattern.

Definition 8 (Overlapping Alphabet Count (OAC)). Let PA be a pattern alphabet and \mathcal{L} be an event log. The overlapping alphabet count of PA in \mathcal{L} is defined as $OAC(PA, \mathcal{L}) = \sum_{\mathbf{t} \in \mathcal{T}} \sum_{p \in [PA]} |\mathcal{I}(p, \mathbf{t})| \cdot f(\mathbf{t})$

Definition 9 (Non-overlapping Alphabet Count (NOAC)). Let PA be a pattern alphabet, and \mathbf{t} be a trace. Let $\mathcal{I}^n(PA, \mathbf{t}) = \{X \subseteq \mathcal{I}(PA, \mathbf{t}) \mid \forall_{i,j \in X} i \geq j + |\text{shortest}([PA], \mathbf{t}, j)|\}$. Let $\mathcal{I}^{n_{max}}(PA, \mathbf{t})$ is an $X \in \mathcal{I}^n(PA, \mathbf{t})$ such that $\nexists X' \in \mathcal{I}^n(PA, \mathbf{t}) X \subset X'$ and $X = \arg \max_{A \in \mathcal{I}^n(PA, \mathbf{t})} |A|$. $\mathcal{I}^{n_{max}}(PA, \mathbf{t})$ denotes the maximal set of all non-overlapping pattern indices of PA in trace \mathbf{t} . The non-overlapping alphabet count of pattern alphabet PA in the event log \mathcal{L} is defined as $NOAC(PA, \mathcal{L}) = \sum_{\mathbf{t} \in \mathcal{T}} |\mathcal{I}^{n_{max}}(PA, \mathbf{t})| \cdot f(\mathbf{t})$

The above definition of NOAC assumes that shorter patterns are preferred over longer patterns. If longer patterns are preferred $\overline{\text{shortest}}([PA], \mathbf{t}, j)$ need to be replaced with $\overline{\text{longest}}([PA], \mathbf{t}, j)$ in the definition.

Definition 10 (Non-overlapping Global Alphabet Count (NOGAC)).

Let \mathbb{PA}_G be the set of all pattern alphabets in the event log and let \mathbf{t} be a trace. Let \mathcal{P} denote the set of all patterns defined by \mathbb{PA}_G i.e., $\mathcal{P} = \bigcup_{PA \in \mathbb{PA}_G} [PA]$. Let

$\mathcal{I}_G = \bigcup_{p \in \mathcal{P}} \mathcal{I}(p, \mathbf{t})$. Let $\mathcal{I}_G^l(\mathcal{P}, \mathbf{t}) = \{X \subset \mathcal{I}_G \mid \forall_{\substack{i \in X \\ j < i \in X}} i > j + |\overline{\text{longest}}(\mathcal{P}, \mathbf{t}, j)|\}$.

$\mathcal{I}_G^{lmax}(\mathcal{P}, \mathbf{t})$ is an $X \in \mathcal{I}_G^l(\mathcal{P}, \mathbf{t})$ such that $\nexists_{X' \in \mathcal{I}_G^l(\mathcal{P}, \mathbf{t})} X \subset X'$ and $X =$

$\arg \max_{A \in \mathcal{I}_G^l(\mathcal{P}, \mathbf{t})} |A|$. Let $\mathcal{I}_G^{nmax}(PA, \mathbf{t}, \mathcal{P}) = \{i \in \mathcal{I}_G^{lmax}(\mathcal{P}, \mathbf{t}) \mid \overline{\text{longest}}(\mathcal{P}, \mathbf{t}, i) \in [PA]\}$. $\mathcal{I}_G^{nmax}(PA, \mathbf{t}, \mathcal{P})$ denotes the maximal set of all non-overlapping pattern indices of PA in trace \mathbf{t} with respect to global alphabet.

The non-overlapping global alphabet count of a pattern alphabet PA in event log \mathcal{L} is defined as $\text{NOGAC}(PA, \mathcal{L}, \mathcal{P}) = \sum_{\mathbf{t} \in \mathcal{L}} |\mathcal{I}_G^{nmax}(PA, \mathbf{t}, \mathcal{P})| \cdot f(\mathbf{t})$

Definition 11 (Conservedness (CON)). Conservedness, of a pattern alphabet PA , measures the degree to which the individual activities involved in the pattern alphabet manifest as the patterns defined by the alphabet and is defined as $\text{CON}(PA, \mathcal{L}) = \frac{\text{NOGAC}(PA, \mathcal{L})}{\mu} \left(1 - \frac{\sigma}{\mu}\right) * 100\%$ where μ and σ denote the mean and standard deviation of the frequency of activities defined by PA in the event log, \mathcal{L} .

The factor $\frac{\sigma}{\mu}$ is the coefficient of variation and measures the dispersion of the distribution of activities in the log. For example, consider three pattern alphabets $\{\mathbf{d}, \mathbf{x}, \mathbf{e}\}$, $\{\mathbf{d}, \mathbf{x}, \mathbf{e}, \mathbf{f}\}$ and $\{\mathbf{d}, \mathbf{x}, \mathbf{e}, \mathbf{h}\}$. Let the equivalence classes of these pattern alphabets be $\{\mathbf{dx}\mathbf{e}\}$, $\{\mathbf{dx}\mathbf{e}\mathbf{f}\}$ and $\{\mathbf{dx}\mathbf{e}\mathbf{h}\}$ respectively. Let the non-overlap alphabet count of these pattern alphabets be $(\{\mathbf{d}, \mathbf{x}, \mathbf{e}\}, 100)$, $(\{\mathbf{d}, \mathbf{x}, \mathbf{e}, \mathbf{f}\}, 60)$, and $(\{\mathbf{d}, \mathbf{x}, \mathbf{e}, \mathbf{h}\}, 40)$. Let the frequency of activities be $(\mathbf{d}, 100)$, $(\mathbf{x}, 100)$, $(\mathbf{e}, 100)$, $(\mathbf{f}, 60)$, and $(\mathbf{h}, 40)$. It is to be noted that not all occurrences of \mathbf{d} , \mathbf{x} , \mathbf{e} manifest as $\mathbf{dx}\mathbf{e}\mathbf{f}$ (or $\mathbf{dx}\mathbf{e}\mathbf{h}$); however, all occurrences of \mathbf{d} , \mathbf{x} , \mathbf{e} manifest as $\mathbf{dx}\mathbf{e}$. Hence, the conservedness value of the pattern alphabets $\{\mathbf{d}, \mathbf{x}, \mathbf{e}\}$, $\{\mathbf{d}, \mathbf{x}, \mathbf{e}, \mathbf{f}\}$, and $\{\mathbf{d}, \mathbf{x}, \mathbf{e}, \mathbf{h}\}$ is 100%, 51% and 30% respectively.

Concurrency in process models affects the manifestation of patterns. For example, let us assume that a sequence involving two activities \mathbf{a} and \mathbf{b} occurs in parallel with another branch. One can notice the pattern \mathbf{ab} manifested in some of the traces in the event log. In other traces, the activities involved in the parallel sub-process might manifest in between \mathbf{a} and \mathbf{b} thus making the NOAC of $\{\mathbf{a}, \mathbf{b}\}$ considerably lower when compared to their mean count. Accordingly conservedness value would be low for this pattern. However, it is to be noted that \mathbf{a} and \mathbf{b} are highly correlated with $\sigma = 0$. In order to compensate for patterns disturbed due to concurrency in process models, we introduce a weighted variant of conservedness, $WCON$, defined as $WCON(PA, \mathcal{L}) = \left(1 - \frac{\sigma}{\mu}\right) * \left(\eta + (1 - \eta) \frac{\text{NOGAC}(PA, \mathcal{L})}{\mu}\right) * 100\%$ where $0 \leq \eta \leq 1$. One can filter

pattern alphabets based on the various metrics presented above. We recommend that pattern alphabets with a high (*weighted*)*conservedness* value ($> 50\%$) be considered as candidates for abstraction.

Though we have defined conservedness using only the *NOAC*, one can define variants by substituting *NOAC* with other count metrics viz., *OAC* and *NOGAC*.

B Annotation of Pattern Graph

The pattern graph can be annotated based on the above pattern metrics for patten selection.

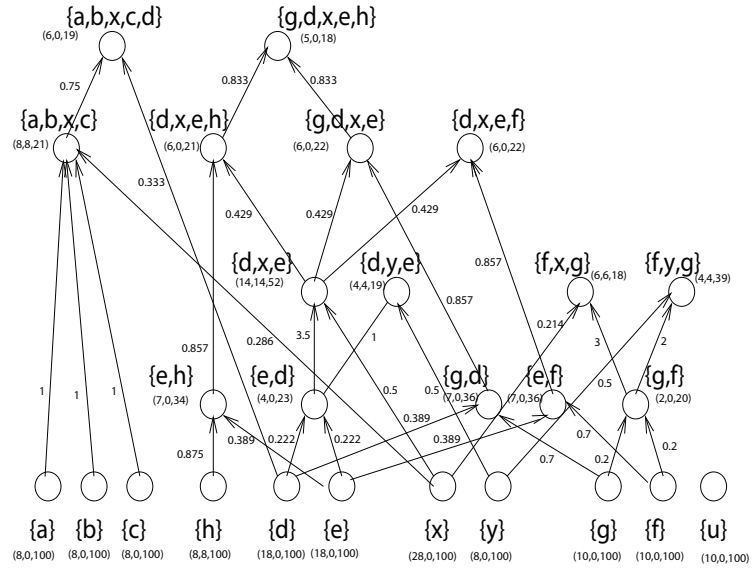


Fig. 6: Pattern graph with node annotation and edge annotation for the example log

- **Node Annotation:** Each pattern node $PA \in V$ is annotated with the triple (*NOAC*, *NOGAC*, *CON*).
- **Edge Annotation:** Each edge $(PA_i, PA_j) \in E$ is annotated with the value $CC(PA_i, PA_j, \mathcal{L})$ where *CC*, the *construction contribution* is defined as $CC(PA_i, PA_j, \mathcal{L}) = NOAC(PA_j, \mathcal{L}) / NOAC(PA_i, \mathcal{L})$. *CC* measures the degree to which the pattern alphabet of a source node of the edge contributes to the pattern alphabet of its target node. The edge annotation is useful in resolving conflicts where a source node is associated to more than one target nodes (i.e. one pattern alphabet can be covered by more than one pattern alphabets.). A source node can be associated to that target node if the edge between them has a larger *CC* value. Figure 6 presents the pattern graph with node annotation and edge annotation for the example log. It is easy

to find that the $CC(\{\mathbf{g}\}, \{\mathbf{g}, \mathbf{d}\}) > CC(\{\mathbf{g}\}, \{\mathbf{g}, \mathbf{f}\})$. Assuming that we use the $\{\mathbf{g}, \mathbf{d}\}$ and $\{\mathbf{g}, \mathbf{f}\}$ in the sub-graph mode and let their abstract activities be \mathbf{X} and \mathbf{Y} respectively, then $\{\mathbf{g}\}$ is assigned to the pattern alphabet $\{\mathbf{g}, \mathbf{d}\}$. With these abstractions, all patterns with pattern alphabets $\{\mathbf{g}\}$, $\{\mathbf{d}\}$ and $\{\mathbf{g}, \mathbf{d}\}$ are represented by \mathbf{X} while the patterns with pattern alphabets $\{\mathbf{f}\}$ and $\{\mathbf{g}, \mathbf{f}\}$ are represented by \mathbf{Y} in all the traces during log transformation.