

# Resolving Uncertain Case Identifiers in Interaction Logs: A User Study<sup>\*</sup>

Marco Pegoraro<sup>✉</sup><sup>[0000-0002-8997-7517]</sup>, Merih Seran Uysal<sup>[0000-0003-1115-6601]</sup>,  
Tom-Hendrik Hülsmann<sup>[0000-0001-8389-5521]</sup>, and Wil M.P. van der  
Aalst<sup>[0000-0002-0955-6940]</sup>

Chair of Process and Data Science (PADS), Department of Computer Science,  
RWTH Aachen, Aachen, Germany  
{pegoraro,uysal,wvdaalst}@pads.rwth-aachen.de  
tom.huelsmann@rwth-aachen.de

**Abstract.** Modern software systems are able to record vast amounts of user actions, stored for later analysis. One of the main types of such user interaction data is click data: the digital trace of the actions of a user through the graphical elements of an application, website or software. While readily available, click data is often missing a case notion: an attribute linking events from user interactions to a specific process instance in the software. In this paper, we propose a neural network-based technique to determine a case notion for click data, thus enabling process mining and other process analysis techniques on user interaction data. We describe our method, show its scalability to datasets of large dimensions, and we validate its efficacy through a user study based on the segmented event log resulting from interaction data of a mobility sharing company. Interviews with domain experts in the company demonstrate that the case notion obtained by our method can lead to actionable process insights.

**Keywords:** Process Mining · Uncertain Event Data · Event-Case Correlation · Case Notion Discovery · Unlabeled Event Logs · Machine Learning · Neural Networks · word2vec · UI Design · UX Design.

## 1 Introduction

In the last decades, the dramatic rise of both performance and portability of computing devices has enabled developers to design software with an ever-increasing level of sophistication. These improvements in computing performance and compactness grew in unison with their access by a larger and larger non-specialized user base, until the point of mass adoption. Such escalation in functionalities caused a subsequent increase in the complexity of software, making it more difficult to access for users. The shift from large screens of desktop computers to

---

<sup>\*</sup> We thank the Alexander von Humboldt (AvH) Stiftung for supporting our research interactions.

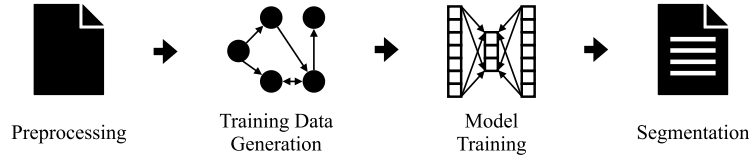
**Table 1.** A sample of click data from the user interactions with the smartphone app of a German mobility sharing company. This dataset is the basis for the qualitative evaluation of the method later presented in this paper (Section 5).

timestamp	screen	user	team	os
2021-01-25 23:00:00.939	pre_booking	b0b00	2070b	iOS
2021-01-25 23:00:03.435	tariffs	b0b00	2070b	iOS
2021-01-25 23:00:04.683	menu	3fc0c	02d1f	Android
2021-01-25 23:00:05.507	my_bookings	3fc0c	02d1f	Android
⋮	⋮	⋮	⋮	⋮

small displays of smartphones, tablets, and other handheld devices has strongly contributed to this increase in the intricacy of software interfaces. *User interface* (UI) design and *user experience* (UX) design aim to address the challenge of managing complexity, to enable users to interact easily and effectively with the software.

In designing and improving user interfaces, important sources of guidance are the records of user interaction data. While in the past enhancement to interfaces were mainly driven by manual intervention of both users of the system and designers, through survey and direct issue reporting in specialized environments, automation in all digital systems have enabled systematic and structured data collection. Many websites and apps track the actions of users, such as pageviews, clicks, and searches. Such type of information is often called *click data*, of which an example is given in Table 1. Click data is a prominent example of *user interaction data*, a digital trace of actions which are recorded, often in real-time, when a user interacts with a system. These can then be analyzed to identify parts of the interface which need to be simplified, through, e.g., frequent item-sets analysis, pattern mining, sequence mining [18], or performance measures such as time spent performing a certain action or visualizing a certain page [13]. However, while such techniques can provide actionable insights with respect to user interface design, they do not account for an important aspect in the system operations: the *process perspective*, a description of all actions in a system contributing to reach a given objective—in the case of user interfaces, the realization of the user’s goal.

A particularly promising sub-field of data science able to account for such perspective of user interfaces is *process mining*. Process mining is a discipline that aims to understand the execution of processes in a data-centric manner, by analyzing collection of historic process executions extracted by information systems—known as *event logs*. Process mining techniques may be used to obtain a model of the process, to measure its conformance with normative behavior, or to analyze the performance of process instances with respect to time and costs. Data from process executions is usually represented as sorted sequences of *events*, each of which is associated with an instance of the process—a *case*. Although the origins of process mining are rooted in the analysis of business process data, in recent years the discipline has been successfully applied to many other contexts,



**Fig. 1.** An overview of the different main phases of our case identifier reconstruction method.

with the goals of obtaining trustworthy descriptive analytics, improving process compliance, increasing time performances, and decreasing costs and wastes. Some examples are logistics [41], auditing [21], production engineering [3], and healthcare [29].

A number of applications of process mining techniques to user interaction data exist—prominently represented by Robotic Process Automation (see Section 6). However, towards the analysis of click data with process mining, a fundamental challenge remains: the association of event data (here, user interactions) with a *process case identifier*. While each interaction logged in a database is associated with a user identifier, which is read from the current active session in the software, there is a lack of an attribute to isolate events corresponding to one single utilization of the software from beginning to end. A function able to subdivide events in sets of single instances of the process, here single utilizations of a software system, is called a *case notion*. Determining the case notion in an event log is a non-trivial task, and is usually a very delicate part of event data extraction from information systems [1]. Aggregating user interactions into cases is of crucial importance, since the case identifier—together with the label of the executed activity and the timestamp of the event—is a fundamental attribute to reconstruct a process instance as a sequence of activities, also known as *control-flow perspective* of a process instance. A vast majority of the process mining techniques available require the control-flow perspective of a process to be known.

In this paper, we propose a novel case attribution approach for click data, an overview of which can be seen in Figure 1. Our method allows us to effectively segment the sequence of interactions from a user into separate cases on the basis of normative behavior. The algorithm takes as input a collection of unsegmented user interaction and the schematic of the system in the form of a link graph, and builds a transition system able to simulate full executions of the process; then, a word2vec neural model is trained on the basis of such simulated full traces, and is then able to split an execution log into well-formed cases. We verify the effectiveness of our method by applying it to a real-life use case scenario related to a mobility sharing smartphone app. Then, we perform common process mining analyses such as process discovery on the resulting segmented log, and we conduct a user study among business owners by presenting the result of such analyses to process experts from the company. Through interviews with such

experts, we assess the impact of process mining analysis techniques enabled by our event-case correlation method. Our evaluation shows that:

- our method obtains a sensible case notion of an input interaction log, using comparatively weak ground truth information;
- our method is efficient, and is able to scale for logs of large dimensions;
- the resulting segmented log provides coherent, actionable insights on the process when analyzed with process mining techniques.

The remainder of the paper is organized as follows. Section 2 presents preliminary concepts and constructs necessary to define our approach. Section 3 illustrates a novel event-case correlation method, which allows to split a stream of interactions into cases—thus enabling process mining analyses on the resulting event log. Section 4 shows the time performance of our method at scale. Section 5 describes the results of our method on a real-life use case scenario related to a mobility sharing app, together with a discussion of interviews of process experts from the company about the impact of process mining techniques enabled by our method. Section 6 examines the current literature, discussing related work and connecting our approach with existing event-case correlation methods. Finally, Section 7 concludes the paper.

## 2 Preliminaries

Let us start by presenting mathematical definitions for the basic structures and concepts necessary for the design of our approach.

### 2.1 Process Mining

*Process mining* is a research field that lies at the intersection of established process sciences such as Business Process Management (BPM) and data science. Its goal is to extract knowledge from so-called *event data* which is continuously collected during the execution of a process. A process can be any sequence of events that are carried out in order to reach a goal. Common examples include business processes such as the *purchase-to-pay* process. However, in recent times, information systems have become ubiquitous and are involved in almost every aspect of modern life. Because of this omnipresence of software systems in processes, they are a prime source for event data. During their execution, such information systems produce large amounts of data in the form of logs that contain information about what actions or tasks were performed at which point in time. Process mining techniques utilize this event data in order to automatically discover new information about the underlying process. This information may then be used in order to improve the observed process in different ways. Despite its young age, the field of process mining already offers a rich ecosystem of algorithms and techniques in areas such as process discovery, conformance checking, process enhancement, and others [2,4].

**Definition 1 (Sequence).** Given a set  $X$ , a finite sequence over  $X$  of length  $n$  is a function  $s \in X^* : \{1, \dots, n\} \rightarrow X$ , and it is written as  $s = \langle s_1, s_2, \dots, s_n \rangle$ . We denote with  $X^*$  the set of all such sequences composed by elements of the set  $X$ . We denote with  $\langle \rangle$  the empty sequence, the sequence with no elements and of length 0. Over the sequence  $s$  we define  $|s| = n$ ,  $s[i] = s_i$  and  $x \in s \Leftrightarrow \exists 1 \leq i \leq n \ s = s_i$ . The concatenation between two sequences is denoted with  $\langle s_1, s_2, \dots, s_n \rangle \cdot \langle s'_1, s'_2, \dots, s'_m \rangle = \langle s_1, s_2, \dots, s_n, s'_1, s'_2, \dots, s'_m \rangle$ . Over the sequence  $\sigma$  of length  $|\sigma| = n$  we define  $hd^k(\sigma) = \langle s_1, \dots, s_{\min(k, n)} \rangle$  to be the function retrieving the first  $k$  elements of the sequence (if possible), and  $tl^k(\sigma) = \langle s_{\max(n-k+1, 1)}, \dots, s_n \rangle$  to be the function retrieving the last  $k$  elements of the sequence (if possible). Note that if  $k \leq 0$  then  $hd^k(\sigma) = tl^k(\sigma) = \langle \rangle$ ; if  $k \geq n$  then  $hd^k(\sigma) = tl^k(\sigma) = \sigma$ ; and for all  $0 \leq k \leq n$  we have that  $hd^k(\sigma) \cdot tl^{n-k}(\sigma) = \sigma$ .

The logs containing the event data that is collected during the execution of the process are called *event logs*. Event logs are a collection of individual events that at least consist of a *timestamp*, the carried out *activity*, and a *case identifier*. These attributes represent the absolute minimum amount of information that is required for most process mining applications. Additionally, there may be other properties associated with the events, for example who carried out the activity or how long its execution did take.

**Definition 2 (Universes).** Let the set  $\mathcal{U}_I$  be the universe of event identifiers. Let the set  $\mathcal{U}_A$  be the universe of activity identifiers. Let the set  $\mathcal{U}_T$  be the totally ordered universe of timestamps. Let the set  $\mathcal{U}_U$  be the universe of users. Let the sets  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_n$  be the universes of attribute domains. The universe of events is defined as  $\mathcal{E} = \mathcal{U}_I \times \mathcal{U}_A \times \mathcal{U}_T \times \mathcal{U}_U \times \mathcal{D}_1 \times \mathcal{D}_2 \times \dots \times \mathcal{D}_n$ .

**Definition 3 (Event and Event Log).** Any element  $e \in \mathcal{E}$  is called an event. Given an event  $e = (i, a, t, u, d_1, \dots, d_n) \in \mathcal{E}$ , we define the following projection functions:  $\pi_I(e) = i$ ,  $\pi_A(e) = a$ ,  $\pi_T(e) = t$ ,  $\pi_U(e) = u$ , and  $\pi_{D_j}(e) = d_j$ . An event log  $L$  is a set  $L \subsetneq \mathcal{E}$  where for any  $e, e' \in L$ , we have  $\pi_I(e) = \pi_I(e') \Rightarrow e = e'$ .

In addition to the events themselves, a case may also be associated metadata that concerns all events of the case and can be used to further describe the underlying process instance (e.g., an order number or a customer identifier).

In order to be able to follow a single process instance throughout the process, each event is normally labeled with a case identifier, an attribute shared among all events belonging to the same process instance—a complete execution of the process to reach a certain objective, specific to each single process. Based on this, the complete event log can be grouped into multiple distinct so-called *cases* that consist of sequences of events with varying lengths. The first event in a case is called the *start event*, while the last event is called the *end event*.

As introduced before, the existence of a timestamp, an activity, and a case identifier is generally a requirement for the majority of process mining operations. Most process mining techniques rely on the fact that a grouping of events

based on the case identifier is possible. For example, consider conformance checking techniques: in order to assess if a process instance is fitting the constraints of the assumed process model, it is a requirement to be able to distinguish between the different process instances. Since this distinction is based on the case identifier, conformance checking is not possible if no such identifier is available. The same is also true for process discovery techniques, in which it is of importance to be able to identify the start and end events. In many areas of application a suitable case identifier is easily available. For example, there might be an order number, a part identifier or a distinct process id. Since these identifiers are in many cases needed during the execution of the process in order to handle the different process instances accordingly, they are generally known to the involved information systems.

However, this is not the case in all circumstances and there exists a significant number of information systems that are involved in processes, but are not process-aware. Examples of such systems include e-mail clients, that may be aware of the recipient but not the concrete case, or machines in production environments that do not have an understanding of the whole production line. In addition to that, there also exist use cases in which the definition of a case is not straightforward and it is therefore not possible to directly assign case identifiers. As introduced before, the analysis of user behavior based on recorded interaction data is an example for such a situation. A case here represents a task that the user performs. At the time of recording, it is not known when a task starts or ends. In such situations, process mining techniques cannot be applied directly to the recorded data. A preprocessing step that correlates events with cases is therefore required.

In contrast to the events in the event log, which model single events in the process, transition systems aim to encode the current state of the process and the transitions between these different states.

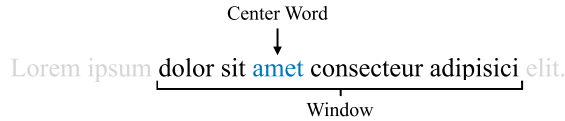
**Definition 4 (Transition System).** *A transition system is a tuple  $TS = (S, A, T, i, S^{end})$  where  $S$  is a set of states that represent a configuration of the process,  $E$  is a set consisting of the actions that can be performed in order to transition between different configurations of the system,  $T \subseteq S \times A \times S$  is a set containing the transitions between configurations,  $i \in S$  is the initial configuration of the process, and  $S^{end} \subseteq S$  is the set of final configurations.*

Starting from the initial state  $i$ , the transition system can move between states according to the transition rules that are defined in  $T$ . A transition system can be obtained from an event log through different types of abstractions. The assumption for these abstractions is that every specific state of the process corresponds to a collection of events in the log. In general, the abstraction is either based on a window of past events, future events, or both. The size of the window is flexible and can be chosen based on the task. When there is more than a single event in the window, one has to additionally choose a representation for the events in the window. Common representations include sets, multisets and sequences of events [5]. Since we will need to quantify the chances of occurring activities, we will attach probabilities to the transitions:

**Definition 5 (Probabilistic Transition System).** A probabilistic transition system is a tuple  $PTS = (S, A, T, i, S^{end})$  where  $S$  is a set of states that represent a configuration of the process,  $A$  is a set consisting of the activities that can be performed in order to transition between different configurations of the process,  $T: S \times A \times S \rightarrow [0, 1]$  is a function expressing the probabilities of transitioning between configurations,  $i \in S$  is the initial configuration of the process, and  $S^{end} \subseteq S$  is the set of final configurations.

## 2.2 Embeddings

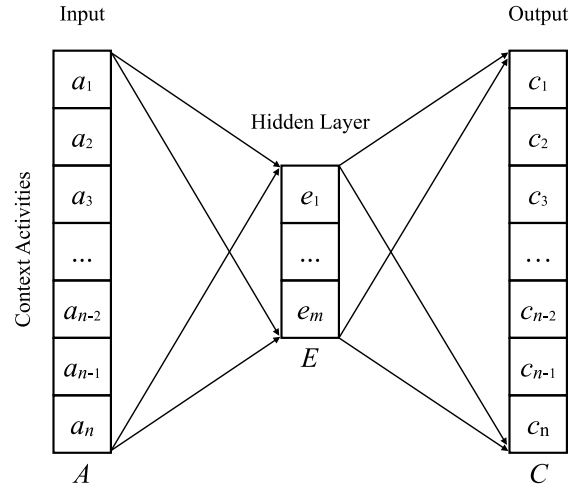
The method presented in this paper is fundamentally based on the concept of *event embeddings* [24], which are themselves based on the natural language processing architecture *word2vec*. The *word2vec* architecture allows the learning of abstract representations of words and their relations, so called *embeddings*. This concept was first proposed in 2013 by Mikolov et. al. in [31] and [32]. The underlying idea of *word2vec* is to encode the relations between words in a body of text using a shallow neural network. The resulting word embeddings are represented by vectors. The more similar the vectors of two words are according to the cosine similarity measure, the more semantically similar the words are. The technique therefore allows to capture the semantic meaning of the words, based on the way they are used in the sentences of a body of text.



**Fig. 2.** An example sentence from a body of text. The window has a size of five and the center word is marked in blue. The two words in front of and after the center word are the context words.

During the training of the two-layer neural network, a sliding window of a specified odd size is used in order to iterate over the sentences. An example for this can be found in Figure 2. The word in the middle of this window is called the *center word*. The words in the window before and after the center word are called *context words*.

There are two different approaches to the *word2vec* architecture; continuous bag-of-words (CBOW) or skip-grams. The main differences between the two approaches are the input and output layers of the network. While in CBOW the frequencies of the context words are used in order to predict the center word, in the skip-gram model the center word is used to predict the context words. The order of the context words is not considered in CBOW. However, the skip-gram model does weigh the context words that are closer to the center word more heavily than those that are further away. A representation of the CBOW architecture can be found in Figure 3.



**Fig. 3.** A graphical representation of the concept behind the event2vec architecture. The vector  $A$  of size  $n$  counts how often every activity occurs in the considered window.  $E$  is the vector representing the event embedding of size  $m$  where  $m \ll n$  and vector  $C$  is a one-hot encoding of the center activity in the ideal case.

Both approaches produce an embedding of the context word in the form of a vector. The advantage of this architecture is that the size of the resulting embedding vectors can be freely determined through the size that is used for the hidden layer. Using this architecture, it is therefore possible to reduce the dimension of the input vector ( $|V|$ ) considerably compared to the output embedding ( $|E|$ ). Additionally, the word embeddings also capture information about the context in which a word is frequently used. As mentioned before, the more similar the vectors of two words, the closer the words are in meaning. In addition to this, the embeddings can also be used in order to predict the center word based on a set of given context words. Because of this versatility, the word2vec architecture is today also widely used in areas other than natural language processing, such as biology [43], medicine [45], or process mining [25].

In the context of process mining, the body of text under consideration is substituted by the event log. In event embeddings, activities and traces take the role of words and sentences in word embeddings. Using this definition, the principle behind word2vec can easily be applied to event data too. Instead of the vocabulary  $V$  there is the set of all possible activities  $A$ . During learning, each activity is associated with its embedding vector  $E$ , which is the output of the hidden layer. The output layer of the network  $C$  ideally represents a one-hot encoding of  $A$ , in which only the desired center activity is mapped to one. Analogous to the word embeddings, event embeddings also capture information about the relations between the different activities. This enables the possibility to find activities that are similar to each other and allows to predict the most likely center activity based on a set of context activities. These properties of event



embeddings are used by the proposed method in order to predict the boundaries between cases, by only using the sequence of activities in the interaction log. As mentioned before, this capability is not only important in the context of process mining, but also in related fields such as robotic process automation which is introduced in more detail in the next section.

### 3 Method

In this section, we illustrate our proposed method for event-case correlation on click data. As mentioned earlier, the goal is to segment the sequence of events corresponding to the interactions of every user in the database into complete process executions (cases). In fact, the click data we consider in this study have a property that we need to account for while designing our method: all events belonging to one case are contiguous in time. Thus, our goal is to determine split points for different cases in a sequence of interactions related to the same user. More concretely, if a user of the app produces the sequence of events  $\langle e_1, e_2, e_3, e_4, e_5, e_6, e_7, e_8, e_9 \rangle$ , our goal is to section such sequence in contiguous subsequences that represent a complete interaction—for instance,  $\langle e_1, e_2, e_3, e_4 \rangle$ ,  $\langle e_5, e_6 \rangle$ , and  $\langle e_7, e_8, e_9 \rangle$ . Such complete interactions should reflect the behavior allowed by the system that supports the process—in the case we examine in our case and user study, such system is a mobile application. We refer to this as the *log segmentation* problem, which can be considered a special case of the event-case correlation problem. In this context, “*unsegmented log*” is synonym with “unlabeled log”.

Rather than being based on a collection of known complete process instances as training set, the creation of our segmentation model is based on behavior described by a model of the system. A type of model particularly suited to the problem of segmentation of user interaction data—and especially click data—is the *link graph*. In fact, since the activities in our process correspond to screens in the app, a graph of the links in the app is relatively easy to obtain, since it can be constructed in an automatic way by following the links between views in the software. This link graph will be the basis for our training data generation procedure.

**Definition 6 (Link Graph).** *A link graph of a software is a graph  $LG = (V, E)$  where  $V$  is the set of pages or screens in the software, and  $E \subseteq V \times V$  represents the links from a page to the next.*

We will use as running example the link graph of Figure 4. The resulting normative traces will then be used to train a neural network model based on the word2vec architecture [32], which will be able to split contiguous user interaction sequences into cases.

#### 3.1 Training Log Generation

To generate the training data, we will begin by exploiting the fact that each process case will only contain events associated with one and only one user.

Let  $L$  be our unsegmented log and  $u \in \mathcal{U}_U$  be a user in  $L$ ; then, we indicate with  $UI_u$  the *user interaction sequence*, a sequence of activities in a sub-log of  $L$  sorted on timestamps where all events are associated with the user  $u$ :  $UI_u = \langle \pi_A(e_1), \pi_A(e_2), \dots, \pi_A(e_n) \rangle$  such that  $e \in UI_u \Rightarrow e \in L \wedge \pi_U(e) = u$ , and it holds<sup>1</sup> that  $\pi_T(e_1) < \pi_T(e_2) < \dots < \pi_T(e_n)$ .

Our training data will be generated by simulating a transition system annotated with probabilities. Initially, for each user  $u \in U$  we create a transition system  $TS_u$  based on the sequence of user interactions  $UI_u$ . The construction of a transition system based on event data is a well-known procedure in process mining [5], which requires to choose a state representation abstraction function  $state: \mathcal{U}_A \rightarrow S_u$  and a window size (or horizon), which are process-specific. In the context of this section, we will show our method using a prefix sequence abstraction with window size 2:  $state(s) = tl^2(s)$ . The application of other abstraction functions is of course possible.

All such transition systems  $TS_u$  share the same initial state  $i$ . To identify the end of sequences, we add a special symbol to the states  $f \in S'$  to which we connect any state  $s \in S$  if it appears at the end of a user interaction sequence. To traverse the transitions to the final state  $f$  we utilize as placeholder the empty label  $\tau$ . Formally, for every user  $u \in \mathcal{U}_U$  and user interaction  $UI_u$  with length  $n = |UI_u|$ , we define  $TS_u = (S_u, A_u, T_u, i_u, S_u^{end})$  as:

- $S_u = \{state(hd^k(UI_u)) \mid 0 \leq k \leq n\} \cup \{f\}$
- $A_u = \{UI_u[k] \mid 0 \leq k \leq n\} \cup \{\tau\}$
- $T_u = \{(state(hd^k(UI_u)), UI_u[k+1], state(hd^{k+1}(UI_u))) \mid 0 \leq k \leq n-1\} \cup \{(state(hd^n(UI_u)), \tau, f)\}$
- $i_u = \langle \rangle$
- $S_u^{end} = \{f\}$

For instance, the user interaction  $\langle M, A, B, C \rangle$  results in  $S_u = \{\langle \rangle, \langle M \rangle, \langle M, A \rangle, \langle A, B \rangle, \langle B, C \rangle, f\}$ ,  $A_u = \{M, A, B, C, \tau\}$ , and  $T_u = \{(\langle \rangle, M, \langle M \rangle), (\langle M \rangle, A, \langle M, A \rangle), (\langle M, A \rangle, B, \langle A, B \rangle), (\langle A, B \rangle, C, \langle B, C \rangle), (\langle B, C \rangle, \tau, f)\}$ .

We then obtain a transition system  $TS' = (S', A', T', i', S'^{end})$  corresponding to the entire log  $L$  by merging the transition systems corresponding to the users:

- $S' = \bigcup_{u \in \mathcal{U}_U} S_u$
- $A' = \bigcup_{u \in \mathcal{U}_U} A_u$
- $T' = \bigcup_{u \in \mathcal{U}_U} T_u$
- $i' = \langle \rangle$
- $S'^{end} = \{f\}$

We also collect information about the frequency of each transition in the log: for the transitions  $(s, a, s') = t \in T$ , we define a weighting function  $\omega: T \rightarrow \mathbb{N}$

<sup>1</sup> We assume user interactions to be tied to clicks on a UI element, so no two user actions can be recorded at the same time. Thus, the total order between events is strict. Assuming a strict total ordering on events is ubiquitous in process mining.

which measures the number of occurrences of the transition  $t$  throughout the entire log:

$$\omega((s, a, s')) = \left| \bigcup_{u \in \mathcal{U}_U} \{(k, u) \mid 0 \leq k \leq n-1 \wedge \text{state}(\text{hd}^k(\text{UI}_u)) = s \wedge \text{UI}_u[k+1] = a \wedge \text{state}(\text{hd}^{k+1}(\text{UI}_u)) = s'\} \right|$$

If  $t \notin T$ ,  $\omega(t) = 0$ . Through  $\omega$ , it is optionally possible to filter out rare behavior by deleting transitions with  $\omega(t) < \epsilon$ , for a small threshold  $\epsilon \in \mathbb{N}$ . Figure 5 shows the transition system  $TS'$  with the chosen abstraction and window size, annotated with both frequencies and transition labels, for the user interactions  $\text{UI}_{u_1} = \langle M, A, M, B, C \rangle$ ,  $\text{UI}_{u_2} = \langle M, B, C, M \rangle$ , and  $\text{UI}_{u_3} = \langle M, A, B, C \rangle$ .

In contrast to transition systems that are created based on logs that are segmented, the obtained transition system might contain states that are not reachable and transitions that are not possible according to the real process. Normally, the transition system abstraction is applied on a case-by-case basis. In our case, however, we applied the abstraction to the whole sequence of interactions that is associated with a specific user, consecutive interactions that belong to different cases will be included as undesired transitions in the transition system. In order to prune undesired transitions from the transition system, we exploit the link graph of the system: a transition in the transition system is only valid if it appears in the link graph. Unreachable states are also pruned.

We will again assume a sequence abstraction. Given a link graph  $LG = (V, E)$ , we define the reduced transition system  $TS_r = (S_r, A_r, T_r, i_r, S_r^{\text{end}})$ , where:

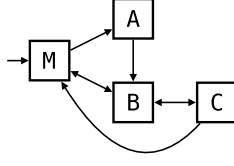
- $S_r = \bigcup_{(s, a, s') \in T} \{s, s'\}$
- $A_r = \{a \in \mathcal{U}_A \mid (s, a, s') \in T'\}$
- $T_r = \{(\langle \dots, a \rangle, a', \langle \dots, a, a' \rangle) \in T' \mid (a, a') \in E\}$
- $i_r = \langle \rangle$
- $S_r^{\text{end}} = \{f\}$

Figure 4 shows a link graph for our running example, and Figure 5 shows how this is used to reduce  $TS'$  into  $TS$ .

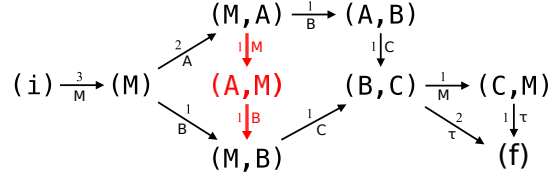
Next, we define probabilities for transitions and states based on the count values for  $\omega(t)$ . Let  $T_{\text{out}} : S \rightarrow \mathcal{P}(T_r)$  be  $T_{\text{out}}(s) = \{(s', a, s'') \in T_r \mid s' = s\}$ ; this function returns all outgoing transitions from a given state. The likelihood of a transition  $(s, a, s') \in T_r$  is then computed with  $l_{\text{trans}} : T_r \rightarrow [0, 1]$ :

$$l_{\text{trans}}(s, a, s') = \frac{\omega(s, a, s')}{\sum_{t \in T_{\text{out}}(s)} \omega(t)}$$

Note that if  $s$  has no outgoing transition and  $T_{\text{out}}(s) = \emptyset$ , we have that  $l_{\text{trans}}(s, a, s') = 0$  for any  $a \in A$  and  $s' \in S_r$ . We will need two more support



**Fig. 4.** The link graph of a simple, fictional system that we are going to use as running example. From this process, we aim to segment the three unsegmented user interactions  $\langle M, A, M, B, C \rangle$ ,  $\langle M, B, C, M \rangle$ , and  $\langle M, A, B, C \rangle$ .



**Fig. 5.** The transition system  $TS'$  obtained by the user interaction data of the example (Figure 4). During the reduction phase, the transition  $(M, A)$  to  $(A, M)$  is removed, since it is not supported by the link graph ( $M$  does not follow  $A$ ). The state  $(A, M)$  is not reachable and is removed entirely (in red). Consequently, the reduced transition system  $TS_r$  is obtained.

functions. We define  $l_{\text{start}}: S_r \rightarrow [0, 1]$  and  $l_{\text{end}}: S_r \rightarrow [0, 1]$  as the probabilities that a state  $s \in S$  is, respectively, the initial and final state of a sequence:

$$l_{\text{start}}(s) = \frac{\sum_{\substack{a \in A \\ s' \in S_r \\ a \in A}} \omega(i, a, s)}{\sum_{\substack{s' \in S_r \\ a \in A}} \omega(s', a, s)} \quad l_{\text{end}}(s) = \frac{\omega(s, \tau, f)}{\sum_{\substack{s' \in S_r \\ a \in A}} \omega(s, a, s')}$$

In our example of Figure 5,  $l_{\text{start}}((M)) = \frac{3}{3} = 1$  and  $l_{\text{end}}((C, M)) = \frac{1}{3}$ .

Such probability functions allow us to define the probabilistic transition system that can simulate an event log based on our dataset of user interactions. We will extend the reduced transition system  $TS_r$  into a probabilistic transition system  $PTS = (S, A, T, i, S^{\text{end}})$  where:

- $S = S_r$
- $A = A_r$
- $T = l_{\text{trans}}$
- $i = i_r$
- $S^{\text{end}} = S_r^{\text{end}}$

Given a path of states  $\langle s_1, s_2, \dots, s_n \rangle$  transitioning in  $PTS$  through the sequence  $\langle (i, a_1, s_1), (s_1, a_2, s_2), \dots, (s_{n-1}, a_n, s_n), (s_n, \tau, f) \rangle$ , we now have the means to compute its probability with the function  $l: S^* \rightarrow [0, 1]$ :

$$l((s_1, s_2, \dots, s_n)) = l_{\text{start}}(s_1) \cdot \prod_{i=2}^n l_{\text{trans}}(s_{i-1}, a_i, s_i) \cdot l_{\text{end}}(s_n)$$

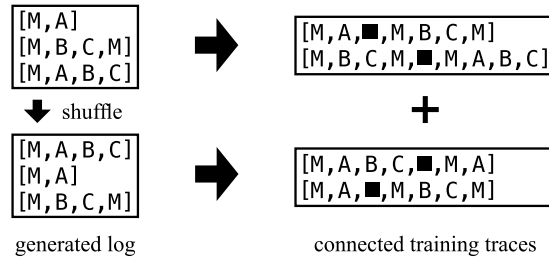
This enables us to obtain an arbitrary number of well-formed process cases as sequences of activities  $\langle a_1, a_2, \dots, a_n \rangle$ , utilizing a Monte Carlo procedure. We can sample a random starting state for the case, through the probability distribution given by  $l_{\text{start}}$ ; then, we compose a path with the probabilities provided by  $l_{\text{trans}}$  and  $l_{\text{end}}$ . The traces sampled in this way will reflect the available user interaction data in terms of initial and final activities, and internal structure, although the procedure still allows for generalization. Such generalization is, however, controlled thanks to the pruning provided by the link graph of the system. We will refer to the set of generated traces as the training log  $L_T$ .

### 3.2 Model Training

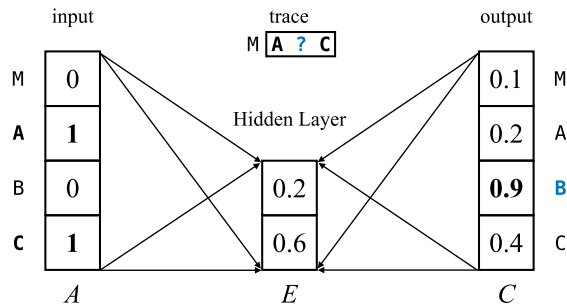
The training log  $L_T$  obtained in Section 3.1 is now used in order to train the segmentation models. The core component of the proposed method consists one or more word2vec models to detect the boundaries between cases in the input log. When applied for natural language processing, the input of a word2vec model is a corpus of sentences which consist of words. Instead of sentences built as sequences of words, we consider traces  $\langle a_1, a_2, \dots, a_n \rangle$  as sequences of activities.

The training log  $L_T$  needs an additional processing step to be used as training set for word2vec. Given two traces  $\sigma_1 \in L_T$  and  $\sigma_2 \in L_T$ , we build a training instance by joining them in a single sequence, concatenating them with a placeholder activity  $\blacksquare$ . So, for instance, the traces  $\sigma_1 = \langle a_1, a_2, a_4, a_5 \rangle \in L_T$  and  $\sigma_2 = \langle a_6, a_7, a_8 \rangle \in L_T$  are combined in the training sample  $\langle a_1, a_2, a_4, a_5, \blacksquare, a_6, a_7, a_8 \rangle$ . This is done repeatedly, shuffling the order of the traces. Figure 6 shows this processing step on the running example.

The word2vec model [32] consists of three layers: an input layer, a single hidden layer, and the output layer. This model has already been successfully employed in process mining to solve the problem of missing events [25]. During training, the network reads the input sequences with a sliding window. The activity occupying the center of the sliding window is called the *center action*, while the surrounding activities are called *context actions*. The proposed method uses the *Continuous Bag-Of-Words* (CBOW) variant of word2vec, where the context actions are introduced as input in the neural network in order to predict the center action. The error measured in the output layer is used for training in order to adjust the weights in the neural network, using the backpropagation algorithm. These forward and backward steps of the training procedure are repeated for all the positions of the sliding window and all the sequences in the training set; when fully trained, the network will output a probability distribution for the center action given the context actions. Figure 7 shows an example of likelihood estimation for a center action in our running example, with a sliding window of size 3.



**Fig. 6.** Construction of the training instances. Traces are shuffled and concatenated with a placeholder end activity.

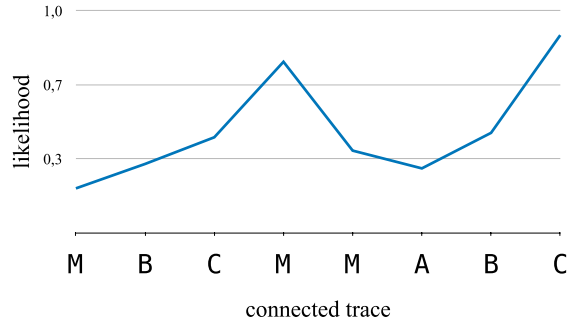


**Fig. 7.** The word2vec neural network. Given the sequence  $\langle A, ?, C \rangle$ , the network produces a probability distribution over the possible activity labels for  $?$ .

### 3.3 Segmentation

Through the word2vec model we trained in Section 3.2, we can now estimate the likelihood of a case boundary ■ at any position of a sequence of user interactions. Figure 8 shows these estimates on one user interaction sequence from the running example. Note that this method of computing likelihoods is easy to extend to an ensemble of predictive models: the different predicted values can be then aggregated, e.g., with the mean or the median.

Next, we use these score to determine case boundaries, which will correspond to prominent peaks in the graph. Let  $\langle p_1, p_2, \dots, p_n \rangle$  be the sequence of likelihoods of a case boundary obtained on a user interaction sequence. We consider  $p_i$  a boundary if it satisfies the following conditions: first,  $p_i > b_1 \cdot p_{i-1}$ ; then,  $p_i > b_2 \cdot p_{i+1}$ ; finally,  $p_i > b_3 \cdot \frac{\sum_{j=i-k-1}^{i-1} p_j}{k}$ , where  $b_1, b_2, b_3 \in [1, \infty)$  and  $k \in \mathbb{N}$  are hyperparameters that influence the sensitivity of the segmentation. The first two inequalities use  $b_1$  and  $b_2$  to ensure that the score is sufficiently higher than the immediate predecessor and successor. The third inequality uses  $b_3$  to make sure that the likelihood is also significantly higher than a neighborhood defined by the parameter  $k$ .



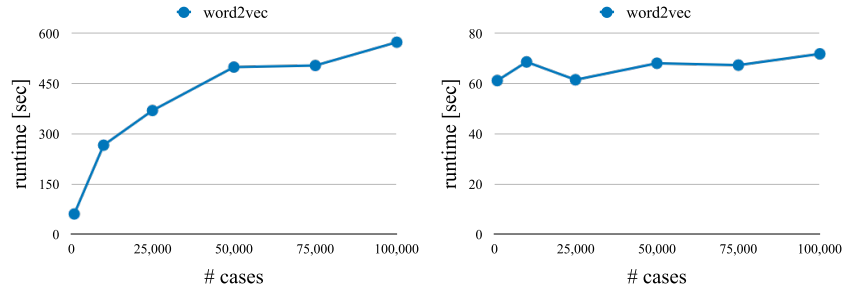
**Fig. 8.** A plot indicating the chances of having a case segment for each position of the user interaction data (second and third trace from the example in Figure 4).

These three conditions allow us to select valid case boundaries within user interaction sequences. Splitting the sequences on such boundaries yields traces of complete process executions, whose events will be assigned a unique case identifier. The set of such traces then constitutes a traditional event log, ready to be analyzed with established process mining techniques.

In the following two sections, we will evaluate two important aspects of our method. Section 4 examines the time performance of the method, and verifies whether it is feasible for large user interaction logs. Section 5 validates our method qualitatively, through a user study in a real-world setting.

## 4 Time performance

Let us now see the efficiency of our method in obtaining a segmentation model. The training phase consists in the generation of the training set and the transition system, and the training of the underlying word2vec models. These steps can take up a considerable amount of time depending on the log size and therefore have to be considered.



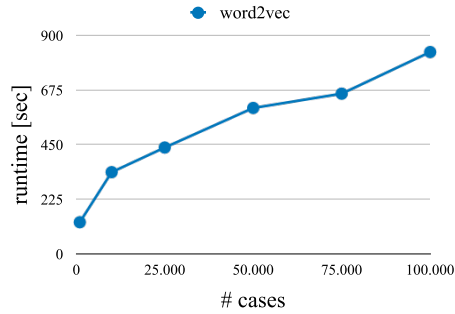
**Fig. 9.** The runtime of the proposed method during the generation of the training log (left) and the time that is required for the model training (right) depending on the number of cases in the input log.

In Figure 9 it can be seen that the time required for the generation of the training set (left) increases quickly for small to medium sized logs, but then plateaus for larger logs. The main factor for the performance of the training set generation is the complexity of the underlying transition system. A larger log will generally contain more behavior, which in turn will lead to a more complex transition system. More paths therefore have to be considered during the generation of the artificial traces. This may explain the plateauing for larger logs; beyond a certain amount of traces, increasing the size of the log will no longer significantly increase the number of variants it contains. The number of states and transitions in the transition system will therefore stop growing, since the system already depicts all of the possible behavior. After this point, the performance of the generation will plateau and is no longer depending on the size of the log.

For the training of the word2vec models, we see a constant required time with minor fluctuations. This indicates that there is no influence of the size of the training log on the performance of the model training. This is caused by the fact that the size of the artificial training log does not depend on the size of the input log, but can be freely chosen. Since the same sized training set was used for all of the logs, the training time did not change significantly.

The combined time that is required for the complete preparation phase of the proposed method, depending on the size of the input log, can be seen in Figure 10. The overall time is mainly influenced by the generation of the transition system, since the model training requires a constant time. Other parts of the





**Fig. 10.** The overall runtime of the proposed method in the preparation phase depending on the number of cases in the input log.

preparation phase such as the computation of the required log statistics have a linear runtime and contribute to the overall runtime behavior that can be seen in Figure 10.

In conclusion, the preparation phase consists of steps with a time complexity of  $\mathcal{O}(T' + S'^2)$  for computing the paths in the underlying transition system  $(S', A, T', i)$  and a constant time complexity (model training). The size of the transition system depends on the size of the input log, but is limited by the number of variants in the log. Overall, it can therefore be said that the time performance of the preparation phase is reasonable (approximately linear in the size of the input) even for larger interaction logs, especially considering that it only has to be performed once, but may be reused for multiple segmentations.

## 5 User Study

In order to validate the utility of process mining workflows in the area of user behavior analysis, a user study was conducted. Such study also aims at assessing the quality of the segmentation produced by the proposed method in a real-life setting, in an area where the ground truth is not available (i.e., there are no normative well-formed cases).

### 5.1 Setting and Methodology

We applied our proposed case segmentation method to a dataset which contains real user interaction data collected from the mobile applications of a German vehicle sharing company. We then utilized the resulting segmented log to analyze user behavior with an array of process mining techniques. Then, the results were presented to process experts from the company, who utilized such results to identify critical areas of the process and suggest improvements. Since the data is from a real-life case study where there is no known ground truth on the actual behavior of the users in the process, we validate our method in a qualitative

way, through an assessment by process experts that the insights obtain through process mining are sensible, truthful, and useful.

In the data, the abstraction for recorded user interactions is the screen (or page) in the app. For each interaction, the system recorded five attributes: `timestamp`, `screen`, `user`, `team`, and `os`. The `timestamp` marks the point in time when the user visited the screen, which is identified by the `screen` attribute, our activity label. The `user` attribute identifies who performed the interaction, and the `team` attribute is an additional field referring to the vehicle provider associated with the interaction. Upon filtering out pre-login screens (not associated with a `user`), the log consists of about 990,000 events originating from about 12,200 users. A snippet of these click data was shown in Table 1, in Section 1.

## 5.2 Results

After applying the segmentation method presented in Section 3 to the click data, as described in the previous section, we analyzed the resulting log with well-known process mining techniques, detailed throughout the section. The findings were presented to and discussed with four experts from the company, consisting of one UX expert, two mobile developers and one manager from a technical area. All of the participants are working directly on the application and are therefore highly familiar with it. We will report here the topics of discussion in the form of questions.

### **Q1: What is the most frequent first screen of an interaction?**

The correct answer to this question is the `station_based_map_dashboard`, which could be computed by considering the first screens for all cases that were identified by the proposed method. All of the participants were able to answer this question correctly. This is expected, as all of the participants are familiar with the application. However, the answers of the participants did not distinguish between the three different types of dashboard that exist in the app. The fact that the map based dashboard is the most frequently used type of dashboard was new and surprising for all of the participants.

### **Q2: What is the most frequent last screen of an interaction?**

The answer to this question can be obtained analogously to that of Q1 directly from the segmented log. In contrast to Q1, not all participants were of the same opinion regarding the answer to this question. Two participants gave the correct answer, which again is the `station_based_map_dashboard`. The other two participants chose the `booking` screen. This screen is the third most frequent case end screen following the `pre_booking` screen. After the correct answer was revealed, one participant proposed that the users may be likely to return to the dashboard after they have completed their goal. This theory can be supported with the available data. It seems that the users have an urge to *clean up* the application and return it to a neutral state before leaving it. Overall, it can be concluded that the participants have a good understanding of the frequent start and end screens of the application. However, the analysis provides more detailed

information and was therefore able to discover aspects about the process that were new for the experts.

**Q3: What is the most frequent interaction with the app?**

This question is asking about the most frequent case variants that are contained in the given log and the associated task of the user. Since the most frequent variants will usually be the shortest variants and a case consisting of only two generic screens cannot be interpreted as a task of the user in a meaningful way, these short variants were not considered for the answer to this question. According to the segmented log, the most common interaction of this type is, selecting a vehicle on the dashboard and checking its availability from the pre-booking screen. One of the four participants did answer this question correctly. Two participants answered that searching for a vehicle on the dashboard is the most frequent interaction, which is closely related to the correct answer but does not include the availability check. The remaining participant answered, opening a booking from the list of all bookings. The results again show that the participants have a good understanding of the usage of the application, but are not able to provide details that are made visible by the log analysis.

**Q4: What is the average length of an interaction with the app?**

For this question, the length of an interaction describes the number of interactions that belong to a case. The correct answer is 4.8 screens, which is rather short. The participants gave the individual answers 50, 30, 12 and 10 screens, which overall results in an average of 25.5. We see that the participants significantly overestimate the length of an average interaction with the app according to the segmented log. However, the average case length is strongly influenced by the employed case attribution method. The mismatch between the results from the log analysis and the expert opinions could therefore be caused by the segmentation that was produced by the proposed method. However, the observed deviations regarding the number of cases were overall not larger than about 50%, which does not explain the large difference between the experts expectations and the calculated value. In order to further examine this, the result was compared to that of a time based segmentation with a fixed threshold of five minutes. These case attribution techniques tend to overestimate the length of cases, as they are not able to distinguish between cases that happen directly after each other. For this reference segmentation, an average case length of 6.7 was calculated. This is comparable to the result of the proposed method and confirms the observation that the experts tend to overestimate the length of interactions significantly.

**Q5: What is the median duration of an interaction with the app?**

For this question, the median duration is used instead of the average, as outliers that have case durations of several days are skewing the average disproportionately. According to the segmented log, the median case duration is 53.4 seconds. The participants gave the answers 240 seconds, 120 seconds, 90 seconds and 60 seconds, leading to an overall average of 127.5 seconds. Similar to the average

length of the interactions, the participants did also overestimate their median duration. Only one participant did give an answer that was close to the calculated value. Both, the significant overestimation of the interaction length and the duration, show that the experts were not able to accurately assess the time a user needs in order to complete a task. This type of analysis is not possible using an unsegmented log and was therefore enabled by the use of the proposed method.

**Q6: How does the median interaction duration on Android and iOS compare?**

As was introduced before, for each interaction it is recorded if it occurred in the Android or iOS application. This allows the comparison between the different applications during analysis. During the analysis it was discovered that the median interaction duration on iOS of 39.4 seconds is significantly shorter than the 92.9 seconds observed for the Android application. The participants were not aware of this difference, as three of the four participants thought that the interaction durations would be the same between the different operating systems and one participant thought that interactions would be shorter on Android. One of the participants argued that Android users may generally be more inclined to “play around” within the application, which may explain the observed difference. Regarding the analysis, the observed deviation could also be caused by differences in the implementation of the screen recording between the two apps. The produced segmentation may reflect cases originating from one of the apps more accurately than those from the other, because the same task of a user may translate to a different sequence of screens in the two apps.

**Q7: Given that 42% of the users use the Android app, what percentage of interactions are from Android users?**

In general one would expect that the fraction of cases that originate from the Android app is similar to the number of users that are using this operating system. The conducted analysis does however show, that only 31% of cases originate from the android app, which is significantly lower than expected. The participants did not expect this uneven distribution, which is emphasized by their answers. Two participants expected a ratio of 50% and two participants answered that 60% of the cases originate from the Android app. In conjunction with the results for the median interaction time that were discussed in Q6/Q7, this means that according to the computed segmentation, Android users tend to use the app longer but overall less frequently.

**Q8: Draw your own process model of the user interactions.**

The participants were asked to draw a *Directly-Follows Graph* (DFG) describing the most common user interactions with the app. A DFG is a simple process model consisting in a graph where activities A and B are connected by an arc if B is executed immediately after A. The concept of this type of graph was explained to the participants beforehand. The experts were given five minutes in

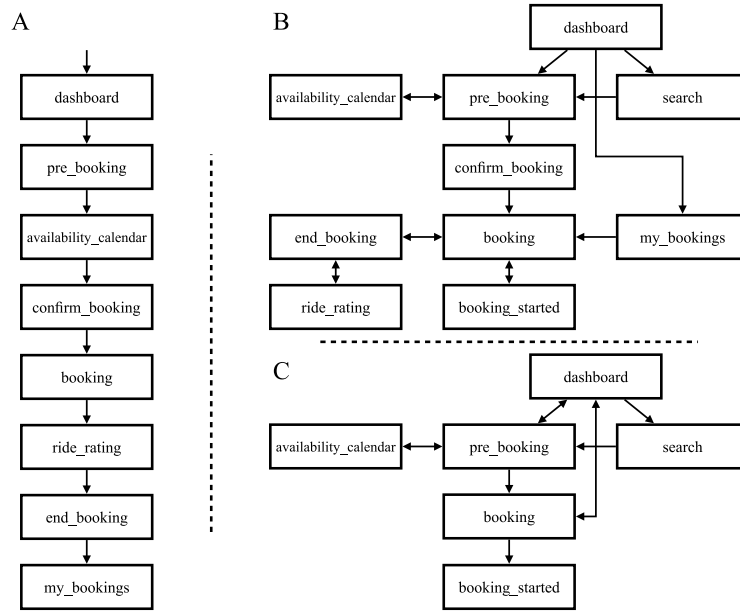


Fig. 11. DFGs created by three of the process experts as part of Q1.

order to create their models. A cleaned up representation of the resulting models can be seen in Figures 11 and 12.

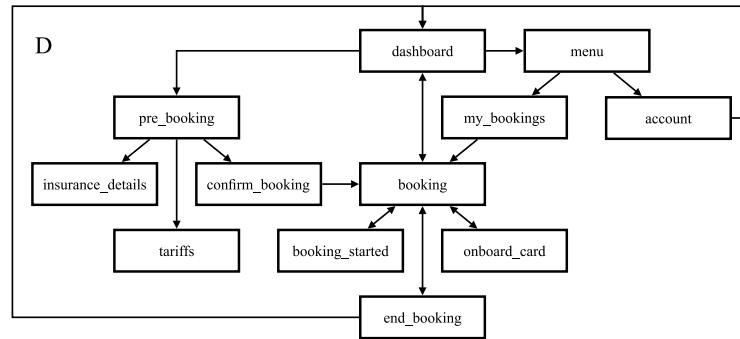
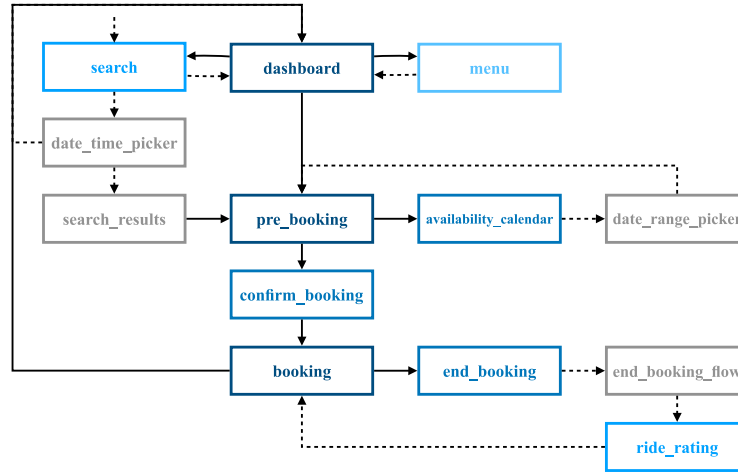


Fig. 12. DFG created by one of the process experts as part of Q1.

For comparison, we created a DFG of the segmented log (Figure 13). Such model was configured to contain a similar amount of different screens as the expert models. The colors indicate the agreement between the model and the expert models. Darker colors signify that a screen was included in more expert



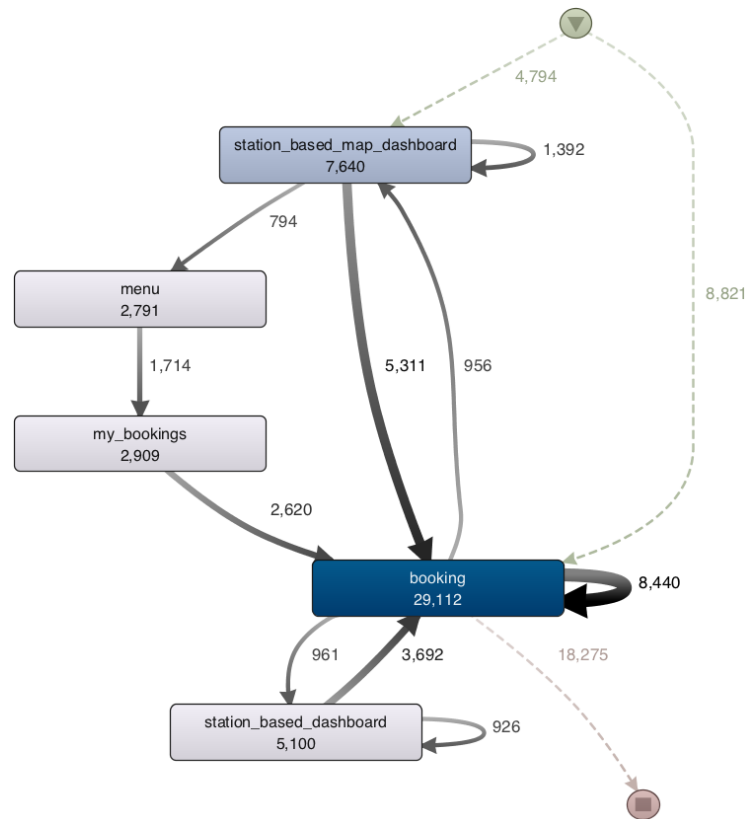
**Fig. 13.** DFG automatically discovered from the log segmented by our method. Darker activities and solid edges were included in models hand-drawn by participants; light-colored activities and dashed edges were not identified by the majority of participants.

models. The dashed edges between the screens signify edges that were identified by the generated model, but are not present in the participant’s models.

The mobile developers (models A and B) tend to describe the interactions in a more precise way that follows the different screens more closely, while the technical manager and UX expert (C and D) provided models that capture the usage of the application in a more abstract way. The fact that the computed model and the expert models are overall very similar to each other suggests that our proposed method is able to create a segmentation that contains cases that are able to accurately describe the real user behavior.

**Q9: Given this process model that is based on interactions ending on the booking screen, what are your observations?**

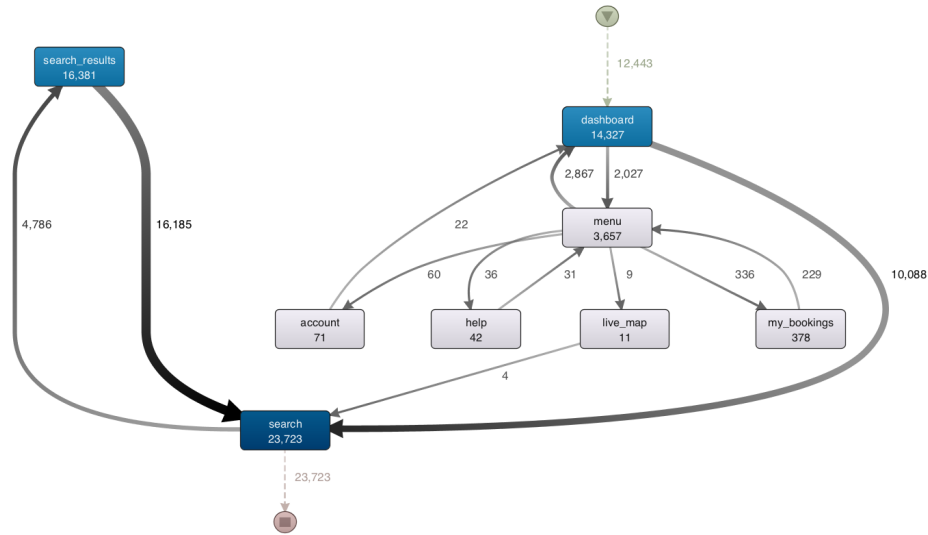
Given the process model shown in Figure 14, the participants were surprised by the fact that the map-based dashboard type is used significantly more frequently than the basic dashboard is surprising to them. Additionally, two of the experts were surprised by the number of users that are accessing their bookings through the list of all bookings (`my_bookings`). This latter observation was also made during the analysis of the segmented log and is the reason that this process model was presented to the experts. In general, a user that has created a booking for a vehicle can access this booking directly from all of the different types of dashboards. The fact that a large fraction of the users takes a detour through the menu and booking list in order to reach the booking screen is therefore surprising. This circumstance was actually already identified by one of the mobile developers some time before this evaluation, while they were manually analyzing the raw interaction recordings data. They noticed this behavior be-



**Fig. 14.** A process model created by using Disco [20], with the **booking** screen as endpoint of the process.

cause they repeatedly encountered the underlying pattern while working with the data for other unrelated reasons. Using the segmented user interaction log, the behavior was however much more discoverable and supported by concrete data rather than just a vague feeling. Another observation that was not made by the participants is that the path through the booking list is more frequently taken by users that originate from the map-based dashboard rather than the basic dashboard. The UX expert suspected that this may have been the case, because the card that can be used to access a booking from the dashboard is significantly smaller on the map-based dashboard and may therefore be missed more frequently by the users. This is a concrete actionable finding of the analysis that was only made possible by the use of process mining techniques in conjunction with the proposed method.

**Q10: Given this process model that is based on interactions ending on the search screen, what are your observations?**



**Fig. 15.** A process model with the `search` screen as endpoint of the process.

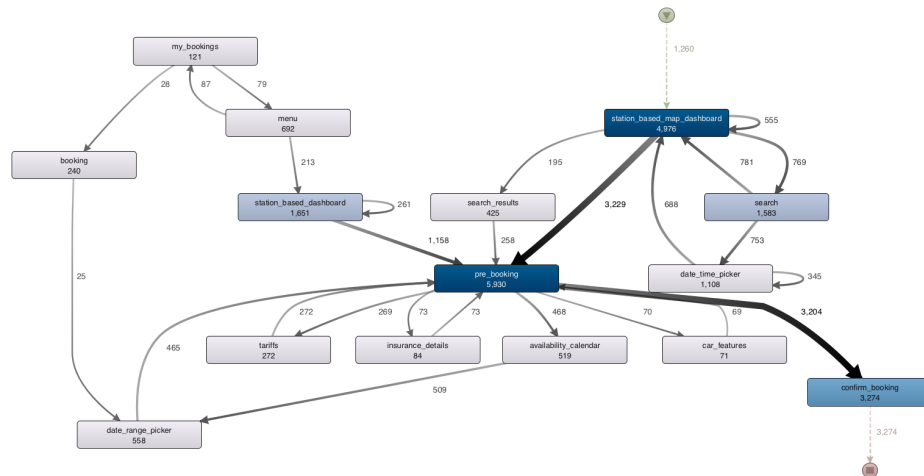
The behavior that was observed during the analysis was tried to be conveyed to the participants using the model that can be found in Figure 15. Since the model is based on all cases including the `search` screen, which start at any type of dashboard, and the `search` screen is directly reachable from the dashboards, it would be expected that no significant amount of other screens are included in the model. This is however not the case, as the `menu` screen and the various screens that are reachable from this screen are included in many of the cases that eventually lead to a search. This suggests that the users that did want to perform a search, tried to find the `search` screen in the main menu, implying that it is not presented prominently enough on the dashboards. None of the experts had this observation when they were presented the discussed model.

#### **Q11: What is the median time a user takes to book a vehicle?**

The correct answer to this question is 66 seconds. This was calculated based on the median time of all cases in which a vehicle booking was confirmed. Three participants gave the answers 420 seconds, 120 seconds and 120 seconds. The fourth participants argued that this time may depend on the type of dashboard that the user is using and answered 300 seconds for the basic dashboard and 120 seconds for the map-based dashboard. When asked to settle on only one time, the participant gave an answer of 180 seconds. Overall this means that the experts estimated a median duration for this task of 3 minutes and 30 seconds. This again is a significant overestimation compared to the value that was obtained by analyzing the real user behavior. Again, a mismatch between the perception of the experts and the real behavior of the users was revealed.



**Q12: Given this process model that is based on interactions ending on the confirm booking screen (Figure 16), what are your observations?**



**Fig. 16.** A process model based on cases that begin in any dashboard and end on the confirm\_booking screen.

Several of the experts observed that the screens that show details about the vehicles and the service, such as `tariffs`, `insurance_details` and `car_features`, are seemingly used much less frequently than expected. In only about 2-10% of cases, the user visits these screens before booking a vehicle. When considering the concrete numbers, the `availability_calendar` screen (which is used to choose a timeframe for the booking) and the `tariffs` screen (which displays pricing information) are used most frequently before a booking confirmation. This suggests that time and pricing information are significantly more important to the users than information about the vehicle or about the included insurance. These findings sparked a detailed discussion between the experts about the possible reasons for the observed behavior. Nonetheless, this shows that models obtained from segmented user interaction logs are an important tool for the analysis of user behavior and that these models provide a valuable foundation for a more detailed analysis by the process experts. Another observation regarding this model was, that a majority of the users seem to choose a vehicle directly from the dashboard cards present on the app rather than using the search functionality. This suggests that the users are more interested in the vehicle itself, rather than looking for any available vehicle at a certain point in time.

**Q13: Discuss the fact that 2% of users activate the intermediate lock before ending the booking.**

The smartphone application offers the functionality to lock certain kinds of vehicles during an active booking. This is for example possible for bicycles, which can be locked by the users during the booking whenever they are leaving the bicycle alone. To do so, the `intermediate_lock` and `intermediate_action` screens are used. During the analysis, it was found that 2% of users use this functionality in order to lock the vehicle directly before ending the booking. This is noteworthy, as it is not necessary to manually lock the vehicle before returning it. All vehicles are automatically locked by the system at the end of each booking. One expert argued that this may introduce additional technical difficulties during the vehicle return, because the system will try to lock the vehicle again. These redundant lock operations, discovered analyzing the segmented log, may introduce errors in the return process.

**Q14: Discuss the fact that only 5% of users visit damages and cleanliness.**

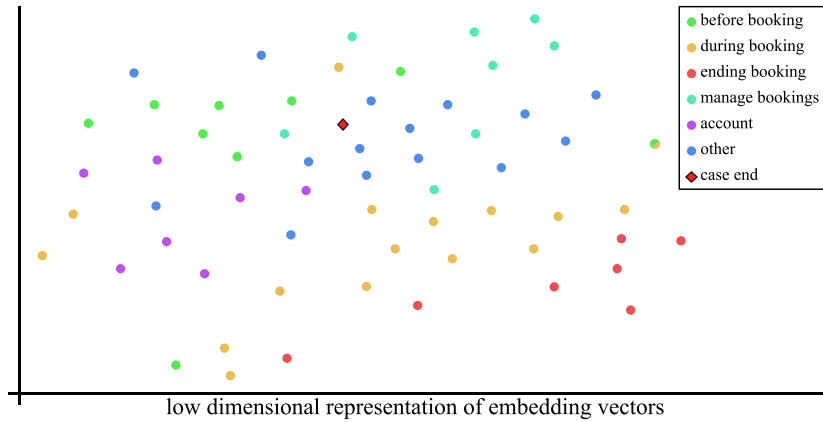
The application allows users to report damages to the vehicles and rate their cleanliness, through the homonymous pages. It was possible to observe that only a small percentage of the users seem to follow this routine, which was surprising to the experts. For the vehicle providers it is generally important that the users are reporting problems with the vehicles; optimally, every user should do this for all of their bookings. According to the data, this is however not the case, as only a small percentage of the users are actually using both of the functionalities. The experts, therefore, concluded that a better communication of these functionalities is required.

### 5.3 Discussion

In this section, we will consider and discuss some aspects, advantages, and limitations of our approach and its applications

In order to evaluate how well the proposed method is able to capture the behavior in the input log and the semantic relationships between activities, we will visualize the embedding vectors of the trained word2vec model. Figure 17 depicts a low dimensional representation of these embedding vectors. The model was trained with the interaction log that was the basis for the conducted case study. The different colors of the dots indicate the different areas of the application. When two actions (dots) are closer to each other in this representation, the actions are related and occur in similar contexts according to the trained model.

Activities that occur during the same phase of the usage will be close to each other in the vector space, and will form clusters. Such clustering of different kinds of actions can be observed in Figure 17. We can see that similar activities indeed form clusters; especially noticeable are the clusters of actions belonging to more distinct phases of the process, such as actions that occur before, during, or at the end of a booking. It can also be observed that the clusters of phases that are more similar to each other are closer to each other in the diagram. For example, the cluster of actions that occur before the booking are closer to those actions that happen during the booking and farther from the ones at the end of the booking. The overall flow of a common interaction with the application is recognizable



**Fig. 17.** A two dimensional representation of the activity embedding vectors of a word2vec model that was trained in the context of the case study. Each dot represents the relative location of an action embedding. The closer two dots are, the more similar are their corresponding embedding vectors. The different colors represent different phases of the process; we can see that similarly colored activities tend to form clusters in the vector space. The dimensional reduction is based on the t-SNE method [28].

in the diagram. This recognizable structure in the activity embedding vectors suggests that the underlying word2vec models is able to abstract the underlying process.

The embedding of the artificial end action that is introduced before model training is marked in Figure 17 with a red rhombus. We can see that it is located near the center of the graphic and shows no clear bias toward any phase of the process. This however also means that end action embedding has no clear relation to any of the clusters. This is expected, as case ends may occur in all of the different phases of the process; however, this can also be considered a weak point in our method, since it indicates that the case end has limited specificity with respect to the type of other activities. One possible solution to this problem that would make the end action more specific is to introduce multiple different end actions, depending on the different process phases, through either different data pre-processing or a post-processing phase on the resulting embeddings.

Even though we applied basic and easily-interpretable process mining techniques to the resulting segmented event log, our user study shows the potential of the application of process mining to user behavior analytics. It was made clear by the study that the process experts are able to comprehend the basic structure of the application and therefore the underlying process well. However, whenever a more detailed view of one aspect of the process was considered, the experts were not able to correctly and accurately assess the real behavior of the users. For instance, concerning the modeling of the process, the experts were able to identify the structure of the most common interactions, but lacked detail

and accuracy. This is especially true when considering the transitions between different screens. The automatically discovered model was more comprehensive and included more behavior and detail.

When the analysis results are processed, visualized and presented to the experts in the right way, they were able to produce clear and actionable results based on the findings. For example it was shown that interactions with the app are much shorter than predicted, that the users are utilizing the bookings list much more frequently than expected, that the map dashboard is the most frequently used dashboard, that the search is less important than the dashboard suggestion cards or that users are unnecessarily locking their vehicles before returning them. Based on these findings, the experts are able to derive concrete and actionable changes to the application, with the goal of improving the overall user experience. Many of the results were completely new and unexpected to the experts and were only enabled through the use of the real dataset in conjunction with the proposed case attribution approach. The time that was required for the segmentation of the large provided interaction log and the subsequent analysis is negligible compared to the amount of information that was obtained.

Overall, the experts were impressed by the findings of the analysis and were able to obtain new insights into the way their users are using the application that were not possible before. Concrete suggestions for improvements could be made and will in the future be implemented in order to improve the user experience of the application, in turn improving the customer satisfaction and lower the required support effort.

## 6 Related Work

### 6.1 Event-Case Correlation

The problem of assigning a case identifier to events in a log is a long-standing challenge in the process mining community [17], and is known by multiple names in literature, including *event-case correlation* problem [10] and *case notion discovery* problem [33]. Event logs where events are missing the case identifier attribute are usually referred to as *unlabeled event logs* [17].

The lack of a case notion has been identified as a major challenge in a number of practical applications, such as analyzing the user interaction with the interface of CT scanners in clinical contexts [42] or measuring the learnability of software systems [30]. Several of the early attempts to solve this problem, such as an early one by Ferreira and Gillblad based on first order Markov models [17], a later approach by Ferreira et al. based on partitioning sequences such that they are minimal and represent a possible process instance [44], or the more recent *Correlation Miner* by Pourmiza et al., based on quadratic programming [39] are very limited in the presence of loops in the process. Other approaches, such as the one by Bayomie et al. [8] can indeed work in the presence of loops, by relying on heuristics based on activities duration which lead to a set of candidate segmented logs. This comes at the cost of a slow computing time. An improvement of the aforementioned method [10] employs simulated annealing to select

an optimal case notion; while still computationally heavy, this method delivers high-quality case attribution results. This was further improved in [9], where the authors reduce the dependence of the method from control flow information and exploit user defined rules to obtain a higher quality result. It is of course important to remember that such methods solve a different and more general problem (the information about the resource is not necessarily available) than the one examined in this paper; in this work, we focus in a more specific setting, where stronger assumptions hold. Such assumptions allow for more efficient segmentation methods, such as the one presented here.

A quite large family of methods approach the problem with a radically different assumption: the hypothesis is that the case information is indeed present in the log, but is hidden. In this context, the case identifier is disguised as a different attribute, or result of a combination of attributes, or learned by applying a similarity function between events. Several such approaches require user-defined rules or domain knowledge to uncover attribute correlations [34,15,16] or require the case notion to be recognizable from a pattern search within the data [6,7].

Many available UI logs are obtained by tracking user action throughout the use of an application, software, or other systems. This means that, similarly to the case study of this paper—which contains roughly one million events—interaction logs are often of large dimensions, at least compared to the typical log sizes in process mining. Therefore, efficiency is important, especially at scale. This motivated our design of a novel method able to reconstruct a case notion for the special case user interaction logs in a fast, interpretable, and loop-robust way, and without relying on ground truth information on cases. This work is an extended version of previous results [38]; we hereby integrate our paper with a more formal description of the method, an evaluation on the time performance of our log segmentation approach, and a full reportage on our mobility app process mining user study.

## 6.2 Uncertain Event Data

The problem of event-case correlation can be positioned in the broader context of *uncertain event data* [35,37]. This research direction aims to analyze event data with imprecise attributes, where single traces might correspond to an array of possible real-life scenarios. For instance, a given event in a log might lack the value of a discrete event attribute such as the activity label, but we might know a set of potential labels; for continuous attributes such as a timestamp, we might have an interval of possible values at our disposal. This type of meta-information on attributes can be quantified with probabilities (*probabilistic uncertainty*) or not (*non-deterministic uncertainty*). Akin to the method proposed in this paper, some techniques allow to obtain probability distributions over such scenarios [36].

Unlabeled logs can then be seen as a specific case of uncertain event logs, where the case identifier is uncertain—since it is not known. Note that having uncertain case identifiers entails more severe consequences than other known types of uncertainty: in all other types, the concept of trace is preserved. According to uncertain event data taxonomies, a missing case identifier can be seen

as a stronger type of *event indetermination* [35], which occurs when an event has been recorded in the log, but it is unclear if it actually happened in reality. Event indetermination is a weaker loss of information than a missing case identifier, in the sense that more information is present and some process mining techniques, albeit specialized, are still possible.

### 6.3 Robotic Process Automation

A notable and rapidly-growing field where the problem of event-case correlation is crucial is *Robotic Process Automation* (RPA), the automation of process activities through the identification of repeated routines in user interactions with software systems [14]. Such routines are automatically discovered from pre-processed user interaction data, then the automatability of such routines is estimated and defined, and software bots are then created to aid the users in repetitive tasks within the process, such as data field completion. As a consequence, the entire discipline of RPA is based on the availability and quality of user interaction logs, which should have a clear and defined case notion. In fact, the problem of case reconstruction is known in the field, and has been identified as a central criticality in automated RPA learning [19] and automated RPA testing [12].

Similarly to many approaches related to the problem at large, existing approaches to event-case correlation in the RPA field often heavily rely on unique start and end events in order to segment the log, either explicitly or implicitly [27,40,26].

### 6.4 Event-Case Correlation Applications

The problem of event-case attribution is different when considered on click data—particularly from mobile apps. Normally, the goal is to learn a function that receives an event as an independent variable and produces a case identifier as an output. In the scenario studied in this paper, however, the user is tracked by the open session in the app during the interaction, and recorded events with different user identifier cannot belong to the same process case. The goal is then to subdivide the sequence of interactions from one user into one or more sessions (cases). While in this user study we assume a prior knowledge of the app where the user interaction is recorded—the link graph—, other ad-hoc techniques to obtain a case notion or segmentation are based on different prior knowledge and different assumptions.

Marrella et al. [30] examined the challenge of obtaining case identifiers for unsegmented user interaction logs in the context of learnability of software systems, by segmenting event sequences with a predefined set of start and end activities as normative information. They find that this approach cannot discover all types of cases, which limits its flexibility and applicability. Jlailaty et al. [23] encounter the segmentation problem in the context of email logs. They segment cases by designing an ad-hoc metric that combines event attributes such as timestamp, sender, and receiver. Their results however show that this method is eluded by

edge cases. Other prominent sources of sequential event data without case attribution are IoT sensors: Janssen et al. [22] address the problem of obtaining process cases from sequential sensor event data by splitting the long traces according to an application-dependent fixed length, to find the optimal sub-trace length such that, after splitting, each case contains only a single activity. One major limitation of this approach that the authors mention is the use of only a single constant length for all of the different activities, which may have varying lengths. More recently, Burattin et al. [11] tackled a segmentation problem for user interactions with a modeling software; in their approach, the segmentation is obtained exploiting eye tracking data, which allows to effectively detect the end of the user interaction with the system.

## 7 Conclusion

In this paper, we showed a case and user study on the topic of the problem of event-case correlation, and presented this problem in the specific application domain of user interaction data.

We examined a case study, the analysis of click data from a mobility sharing smartphone application. To perform log segmentation, we proposed an original technique based on the word2vec neural network architecture, which can obtain case identification for an unlabeled user interaction log on the sole basis of a link graph of the system as normative information. We then presented a user study, where experts of the process were confronted with insights obtained by applying process mining techniques to the log segmented using our method. The interviews with experts confirm that our technique helped to uncover hidden characteristics of the process, including inefficiencies and anomalies unknown to the domain knowledge of the business owners. Importantly, the analyses yielded actionable suggestions for UI/UX improvements, some of which were readily incorporated in the mobile app. This substantiates the scientific value of event-log correlation techniques for user interaction data, and shows the direct benefits of the application of process analysis techniques to data from the user interaction domain. Furthermore, the user study demonstrates the validity of the segmentation method presented in this paper, and its ability of producing a coherent case notion via the segmentation of user interaction sequences. Quantitative experiments with logs of increasing size show the scalability of our method, which is able to preserve its time performance with logs of large dimensions. Lastly, we highlighted how the use of a word2vec model results in a fixed-length representation for activities which expresses some of the semantic relationships between the respective activity labels.

As future work, we intend to further validate our technique by lifting it from the scope of a user study by means of a quantitative evaluation on its efficacy, to complement the qualitative one showed in this paper. Since our segmentation technique has several points of improvement, including the relatively high number of hyperparameters, it would benefit from a heuristic procedure to determine the (starting) value for such hyperparameters. It is also possible to apply differ-

ent encoding techniques for embeddings in place of word2vec, which may results in a better segmentation quality for specific interaction logs. Finally, other future work may consider additional event data perspectives, such as adding the data perspective to our technique by encoding additional attributes in the training set of the neural network model.

## References

1. van der Aalst, W.M.P.: *Process Mining - Data Science in Action*, Second Edition. Springer (2016). <https://doi.org/10.1007/978-3-662-49851-4>
2. van der Aalst, W.M.P., Adriansyah, A., de Medeiros, A.K.A., et al.: *Process mining manifesto*. In: *Business Process Management Workshops - BPM 2011 International Workshops*, Clermont-Ferrand, France, August 29, 2011, Revised Selected Papers, Part I. *Lecture Notes in Business Information Processing*, vol. 99, pp. 169–194. Springer (2011). [https://doi.org/10.1007/978-3-642-28108-2\\_19](https://doi.org/10.1007/978-3-642-28108-2_19)
3. van der Aalst, W.M.P., Brockhoff, T., Ghahfarokhi, A.F., Pourbafrani, M., Uysal, M.S., van Zelst, S.J.: *Removing operational friction using process mining: Challenges provided by the internet of production (iop)*. In: Hammoudi, S., Quix, C., Bernardino, J. (eds.) *Data Management Technologies and Applications - 9th International Conference, DATA 2020, Virtual Event, July 7-9, 2020, Revised Selected Papers. Communications in Computer and Information Science*, vol. 1446, pp. 1–31. Springer (2020). [https://doi.org/10.1007/978-3-030-83014-4\\_1](https://doi.org/10.1007/978-3-030-83014-4_1)
4. van der Aalst, W.M.P., Carmona, J. (eds.): *Process Mining Handbook, Lecture Notes in Business Information Processing*, vol. 448. Springer (2022). <https://doi.org/10.1007/978-3-031-08848-3>
5. van der Aalst, W.M.P., Rubin, V.A., Verbeek, H.M.W., van Dongen, B.F., Kindler, E., Günther, C.W.: *Process mining: a two-step approach to balance between underfitting and overfitting*. *Software and Systems Modeling* **9**(1), 87–111 (2010). <https://doi.org/10.1007/s10270-008-0106-z>
6. Andaloussi, A.A., Burattin, A., Weber, B.: *Toward an automated labeling of event log attributes*. In: *Enterprise, Business-Process and Information Systems Modeling - 19th International Conference, BPMDS 2018, 23rd International Conference, EMMSAD 2018, Held at CAiSE 2018, Tallinn, Estonia, June 11-12, 2018, Proceedings. Lecture Notes in Business Information Processing*, vol. 318, pp. 82–96. Springer (2018). [https://doi.org/10.1007/978-3-319-91704-7\\_6](https://doi.org/10.1007/978-3-319-91704-7_6)
7. Bala, S., Mendling, J., Schimak, M., Queteschiner, P.: *Case and activity identification for mining process models from middleware*. In: *The Practice of Enterprise Modeling - 11th IFIP WG 8.1. Working Conference, PoEM 2018, Vienna, Austria, October 31 - November 2, 2018, Proceedings. Lecture Notes in Business Information Processing*, vol. 335, pp. 86–102. Springer (2018). [https://doi.org/10.1007/978-3-030-02302-7\\_6](https://doi.org/10.1007/978-3-030-02302-7_6)
8. Bayomie, D., Awad, A., Ezat, E.: *Correlating unlabeled events from cyclic business processes execution*. In: *Advanced Information Systems Engineering - 28th International Conference, CAiSE 2016, June 13-17, 2016. Proceedings. Lecture Notes in Computer Science*, vol. 9694, pp. 274–289. Springer (2016). [https://doi.org/10.1007/978-3-319-39696-5\\_17](https://doi.org/10.1007/978-3-319-39696-5_17)
9. Bayomie, D., Ciccio, C.D., Mendling, J.: *Event-case correlation for process mining using probabilistic optimization*. *CoRR* **abs/2206.10009** (2022). <https://doi.org/10.48550/arXiv.2206.10009>



10. Bayomie, D., Ciccio, C.D., Rosa, M.L., Mendling, J.: A probabilistic approach to event-case correlation for process mining. In: Conceptual Modeling - 38th International Conference, ER 2019, November 4-7, 2019, Proceedings. Lecture Notes in Computer Science, vol. 11788, pp. 136–152. Springer (2019). [https://doi.org/10.1007/978-3-030-33223-5\\_12](https://doi.org/10.1007/978-3-030-33223-5_12)
11. Burattin, A., Kaiser, M., Neurauter, M., Weber, B.: Learning process modeling phases from modeling interactions and eye tracking data. *Data & Knowledge Engineering* **121**, 1–17 (2019). <https://doi.org/10.1016/j.datak.2019.04.001>
12. Chacón-Montero, J., Jiménez-Ramírez, A., Enríquez, J.G.: Towards a method for automated testing in robotic process automation projects. In: Proceedings of the 14th International Workshop on Automation of Software Test, AST@ICSE 2019, May 27, 2019, Montreal, QC, Canada. pp. 42–47. IEEE / ACM (2019). <https://doi.org/10.1109/AST.2019.00012>
13. Dhandi, M., Chakrawarti, R.K.: A comprehensive study of web usage mining. In: 2016 Symposium on Colossal Data Analysis and Networking (CDAN). pp. 1–5. IEEE (2016). <https://doi.org/10.1109/CDAN.2016.7570889>
14. Dumas, M., Rosa, M.L., Leno, V., Polyvyanyy, A., Maggi, F.M.: Robotic process mining. In: Process Mining Handbook, Lecture Notes in Business Information Processing, vol. 448, pp. 468–491. Springer (2022). [https://doi.org/10.1007/978-3-031-08848-3\\_16](https://doi.org/10.1007/978-3-031-08848-3_16)
15. Engel, R., Bose, R.P.J.C.: A case study on analyzing inter-organizational business processes from EDI messages using physical activity mining. In: 47th Hawaii International Conference on System Sciences, HICSS 2014, Waikoloa, HI, USA, January 6-9, 2014. pp. 3858–3867. IEEE Computer Society (2014). <https://doi.org/10.1109/HICSS.2014.479>
16. Engel, R., Krathu, W., Zapletal, M., Pichler, C., Bose, R.P.J.C., van der Aalst, W.M.P., Werthner, H., Huemer, C.: Analyzing inter-organizational business processes - process mining and business performance analysis using electronic data interchange messages. *Information Systems and e-Business Management* **14**(3), 577–612 (2016). <https://doi.org/10.1007/s10257-015-0295-2>
17. Ferreira, D.R., Gillblad, D.: Discovering process models from unlabelled event logs. In: Business Process Management, 7th International Conference, BPM 2009, September 8-10, 2009. Proceedings. Lecture Notes in Computer Science, vol. 5701, pp. 143–158. Springer (2009). [https://doi.org/10.1007/978-3-642-03848-8\\_11](https://doi.org/10.1007/978-3-642-03848-8_11)
18. Fournier-Viger, P., Wu, C., Tseng, V.S., Cao, L., Nkambou, R.: Mining partially-ordered sequential rules common to multiple sequences. *IEEE Transactions on Knowledge and Data Engineering* **27**(8), 2203–2216 (2015). <https://doi.org/10.1109/TKDE.2015.2405509>
19. Gao, J., van Zelst, S.J., Lu, X., van der Aalst, W.M.P.: Automated robotic process automation: A self-learning approach. In: On the Move to Meaningful Internet Systems: OTM 2019 Conferences - Confederated International Conferences: CoopIS, ODBASE, C&TC 2019, Rhodes, Greece, October 21-25, 2019, Proceedings. Lecture Notes in Computer Science, vol. 11877, pp. 95–112. Springer (2019). [https://doi.org/10.1007/978-3-030-33246-4\\_6](https://doi.org/10.1007/978-3-030-33246-4_6)
20. Günther, C.W., Rozinat, A.: Disco: Discover your processes. In: Proceedings of the Demonstration Track of the 10th International Conference on Business Process Management (BPM 2012), Tallinn, Estonia, September 4, 2012. CEUR Workshop Proceedings, vol. 940, pp. 40–44. CEUR-WS.org (2012), <http://ceur-ws.org/Vol-940/paper8.pdf>

21. Jans, M., Eulerich, M.: Process mining for financial auditing. In: van der Aalst, W.M.P., Carmona, J. (eds.) *Process Mining Handbook, Lecture Notes in Business Information Processing*, vol. 448, pp. 445–467. Springer (2022). [https://doi.org/10.1007/978-3-031-08848-3\\_15](https://doi.org/10.1007/978-3-031-08848-3_15)
22. Janssen, D., Mannhardt, F., Koschmider, A., van Zelst, S.J.: Process model discovery from sensor event data. In: *Process Mining Workshops - ICPM 2020 International Workshops*, October 5-8, 2020, Revised Selected Papers. *Lecture Notes in Business Information Processing*, vol. 406, pp. 69–81. Springer (2020). [https://doi.org/10.1007/978-3-030-72693-5\\_6](https://doi.org/10.1007/978-3-030-72693-5_6)
23. Jlalaty, D., Grigori, D., Belhajjame, K.: Business process instances discovery from email logs. In: *2017 IEEE International Conference on Services Computing, SCC 2017*, June 25-30, 2017. pp. 19–26. IEEE Computer Society (2017). <https://doi.org/10.1109/SCC.2017.12>
24. Koninck, P.D., vanden Broucke, S., Weerdt, J.D.: act2vec, trace2vec, log2vec, and model2vec: Representation learning for business processes. In: *Business Process Management - 16th International Conference, BPM 2018, Sydney, NSW, Australia, September 9-14, 2018, Proceedings. Lecture Notes in Computer Science*, vol. 11080, pp. 305–321. Springer (2018). [https://doi.org/10.1007/978-3-319-98648-7\\_18](https://doi.org/10.1007/978-3-319-98648-7_18)
25. Lakhani, K., Narayan, A.: A neural word embedding approach to system trace reconstruction. In: *2019 IEEE International Conference on Systems, Man and Cybernetics, SMC*, October 6-9, 2019. pp. 285–291. IEEE (2019). <https://doi.org/10.1109/SMC.2019.8914322>
26. Leno, V., Augusto, A., Dumas, M., Rosa, M.L., Maggi, F.M., Polyvyanyy, A.: Identifying candidate routines for robotic process automation from unsegmented UI logs. In: *2nd International Conference on Process Mining, ICPM 2020*, October 4-9, 2020. pp. 153–160. IEEE (2020). <https://doi.org/10.1109/ICPM49681.2020.00031>
27. Linn, C., Zimmermann, P., Werth, D.: Desktop activity mining - A new level of detail in mining business processes. In: *48. Jahrestagung der Gesellschaft für Informatik, Architekturen, Prozesse, Sicherheit und Nachhaltigkeit, INFORMATIK 2018 - Workshops*, September 26-27, 2018. *LNI*, vol. P-285, pp. 245–258. GI (2018), <https://dl.gi.de/20.500.12116/17225>
28. van der Maaten, L., Hinton, G.: Visualizing data using t-SNE. *Journal of Machine Learning Research* **9**(11) (2008), <http://jmlr.org/papers/v9/vandermaaten08a.html>
29. Mans, R., van der Aalst, W.M.P., Vanwersch, R.J.B.: *Process Mining in Healthcare - Evaluating and Exploiting Operational Healthcare Processes*. Springer Briefs in Business Process Management, Springer (2015). <https://doi.org/10.1007/978-3-319-16071-9>
30. Marrella, A., Catarci, T.: Measuring the learnability of interactive systems using a Petri Net based approach. In: *Proceedings of the 2018 on Designing Interactive Systems Conference, DIS*, June 09-13, 2018. pp. 1309–1319. ACM (2018). <https://doi.org/10.1145/3196709.3196744>
31. Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. In: *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings* (2013), <http://arxiv.org/abs/1301.3781>
32. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems. Proceedings of a meeting held De-*

- cember 5-8, 2013 (2013), <https://proceedings.neurips.cc/paper/2013/hash/9aa42b31882ec039965f3c4923ce901b-Abstract.html>
33. de Murillas, E.G.L., Reijers, H.A., van der Aalst, W.M.P.: Case notion discovery and recommendation: automated event log building on databases. *Knowledge and Information Systems* **62**(7), 2539–2575 (2020). <https://doi.org/10.1007/s10115-019-01430-6>
  34. Nezhad, H.R.M., Saint-Paul, R., Casati, F., Benatallah, B.: Event correlation for process discovery from web service interaction logs. *VLDB J.* **20**(3), 417–444 (2011). <https://doi.org/10.1007/s00778-010-0203-9>
  35. Pegoraro, M.: Process mining on uncertain event data. In: International Conference on Process Mining ICPM 2021, Doctoral Consortium and Tool Demonstration Track, Eindhoven, the Netherlands, October 31–November 4, 2021. CEUR Workshop Proceedings, vol. 3098, pp. 1–2. CEUR-WS.org (2021), [http://ceur-ws.org/Vol-3098/dc\\_176.pdf](http://ceur-ws.org/Vol-3098/dc_176.pdf)
  36. Pegoraro, M., Bakullari, B., Uysal, M.S., van der Aalst, W.M.P.: Probability estimation of uncertain process trace realizations. In: Process Mining Workshops - ICPM 2021 International Workshops, October 31 - November 4, 2021, Revised Selected Papers. Lecture Notes in Business Information Processing, vol. 433, pp. 21–33. Springer (2021). [https://doi.org/10.1007/978-3-030-98581-3\\_2](https://doi.org/10.1007/978-3-030-98581-3_2)
  37. Pegoraro, M., Uysal, M.S., van der Aalst, W.M.P.: PROVED: A tool for graph representation and analysis of uncertain event data. In: Application and Theory of Petri Nets and Concurrency - 42nd International Conference, PETRI NETS 2021, June 23-25, 2021, Proceedings. Lecture Notes in Computer Science, vol. 12734, pp. 476–486. Springer (2021). [https://doi.org/10.1007/978-3-030-76983-3\\_24](https://doi.org/10.1007/978-3-030-76983-3_24)
  38. Pegoraro, M., Uysal, M.S., Hülsmann, T., van der Aalst, W.M.P.: Uncertain case identifiers in process mining: A user study of the event-case correlation problem on click data. In: Enterprise, Business-Process and Information Systems Modeling - 23rd International Conference, BPMDS 2022 and 27th International Conference, EMMSAD 2022, Held at CAiSE 2022, Leuven, Belgium, June 6-7, 2022, Proceedings. Lecture Notes in Business Information Processing, vol. 450, pp. 173–187. Springer (2022). [https://doi.org/10.1007/978-3-031-07475-2\\_12](https://doi.org/10.1007/978-3-031-07475-2_12)
  39. Pourmirza, S., Dijkman, R.M., Grefen, P.: Correlation miner: Mining business process models and event correlations without case identifiers. *International Journal of Cooperative Information Systems* **26**(2), 1742002:1–1742002:32 (2017). <https://doi.org/10.1142/S0218843017420023>
  40. Ramirez, A.J., Reijers, H.A., Barba, I., Valle, C.D.: A method to improve the early stages of the robotic process automation lifecycle. In: Advanced Information Systems Engineering - 31st International Conference, CAiSE, June 3-7, 2019, Proceedings. Lecture Notes in Computer Science, vol. 11483. Springer (2019). [https://doi.org/10.1007/978-3-030-21290-2\\_28](https://doi.org/10.1007/978-3-030-21290-2_28)
  41. Reikemeyer, L. (ed.): *Process Mining in Action: Principles, Use Cases and Outlook*. Springer Cham (2020). <https://doi.org/10.1007/978-3-030-40172-6>
  42. Reindler, J.: Siemens healthineers: Process mining as an innovation driver in product management. In: *Process Mining in Action*, pp. 143–157. Springer (2020)
  43. Tsukiyama, S., Hasan, M.M., Fujii, S., Kurata, H.: LSTM-PHV: prediction of human-virus protein-protein interactions by LSTM with word2vec. *Briefings Bioinform.* **22**(6) (2021). <https://doi.org/10.1093/bib/bbab228>
  44. Walicki, M., Ferreira, D.R.: Sequence partitioning for process mining with unlabeled event logs. *Data & Knowledge Engineering* **70**(10), 821–841 (2011). <https://doi.org/10.1016/j.datak.2011.05.003>

45. Zhang, Y.F., Wang, X., Kaushik, A.C., Chu, Y., Shan, X., Zhao, M.Z., Xu, Q., Wei, D.Q.: SPVec: A word2vec-inspired feature representation method for drug-target interaction prediction. *Frontiers in Chemistry* **7** (2020). <https://doi.org/10.3389/fchem.2019.00895>, <https://www.frontiersin.org/articles/10.3389/fchem.2019.00895>