

The Impact of Event Data Partitioning on Privacy-aware Process Discovery

Jungeun Lim¹, Stephan A. Fahrenkrog-Petersen^{2,3}, Xixi Lu⁴, Jan Mendling^{2,3},
and Minseok Song¹

¹ Pohang University of Science and Technology, Pohang, South Korea

² Weizenbaum Institute, Berlin, Germany

³ Humboldt-Universität zu Berlin, Berlin, Germany

⁴ Utrecht University, Utrecht, The Netherlands

stephan.fahrenkrog-petersen@hu-berlin.de

Abstract. Information systems support the execution of business processes. The event logs of these executions generally contain sensitive information about customers, patients, and employees. The corresponding privacy challenges can be addressed by anonymizing the event logs while still retaining *utility* for process discovery. However, trading off utility and privacy is difficult: the higher the complexity of event log, the higher the loss of utility by anonymization. In this work, we propose a pipeline that combines anonymization and *event data partitioning*, where event abstraction is utilized for partitioning. By leveraging event abstraction, event logs can be segmented into multiple parts, allowing each sub-log to be anonymized separately. This pipeline preserves privacy while mitigating the loss of utility. To validate our approach, we study the impact of event partitioning on two anonymization techniques using three real-world event logs and two process discovery techniques. Our results demonstrate that event partitioning can bring improvements in process discovery utility for directly-follows-based anonymization techniques.

Keywords: Privacy-preserving Process Mining · Event Log Anonymization · Event Log Pre-processing

1 Introduction

Event logs recorded by information systems are the backbone of process mining. Typically, these event logs contain personal information [32] that need to be protected in terms of privacy [10]. Such protection can be achieved through the anonymization of event logs [12]. Anonymization aims to provide a formal privacy guarantee for the data and, at the same time, maximize the *utility* of the anonymized event log. One way utility can be defined is by evaluating the quality of process models discovered using the anonymized log [27].

Many anonymization techniques work through the insertion of noise [14] with the aim of providing *differential privacy* [6]. This noise is inserted into the result of queries that determine how often a control-flow behavior appears in the original event log. The query can either return a count of the trace-variants or the

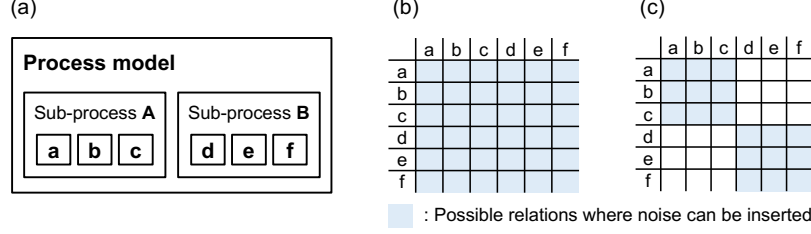


Fig. 1. Noise insertion differences between considering and not considering process structure. (a): Structure of process model, (b): Noise insertion areas when not considering process structure, (c): Noise insertion areas when considering process structure.

directly-follows relations within an event log [24]. A common goal of anonymization is to minimize the inserted noise while maximizing the preserved utility. A high-level of utility is especially hard to reach for less-structured processes.

In this paper, we build on the hypothesis that event data partitioning could facilitate a better trade-off between privacy and utility by reducing the amount of required noise. As an example, consider a process where sub-process *A* consists of activities *a*, *b*, and *c*, and sub-process *B* consists of activities *d*, *e*, and *f*, as shown in Figure 1 (a). If a privacy-preserving technique injects noise into the directly-follows relations without considering the process structure, the resulting noise can be extensive, as illustrated in Figure 1 (b) in blue. However, when the structure is considered, as in Figure 1 (c) in blue, the noise has only to be applied to each sub-process independently. The additional noise in Figure 1 (b) could introduce unrealistic process behaviors, impacting process discovery utility.

In line with this argument, we propose a pipeline that first partitions an event log and then applies anonymization. The anonymized logs are subsequently used for process discovery. The utility of the anonymized log is assessed based on the quality of the discovered models. First, we assess its effectiveness by measuring the utility of process models derived from sub-logs. We compare the utility of these models with those obtained from only anonymized event logs. Second, we evaluate whether event data partitioning always provides utility benefits. To do so, we compare the utility of process models obtained by applying partitioning before anonymization versus those obtained by applying partitioning after anonymization. We evaluated these questions for different types of anonymization, considering both directly-follows-based and trace-variant-based anonymization. Our results show that performing event partitioning before anonymization leads to better process discovery utility for directly-follows-based anonymization.

The remainder of this paper is structured as follows: Section 2 reviews existing research on privacy-preserving process mining and event log abstraction. Section 3 introduces our proposed pipeline, and Section 4 presents the evaluation. Finally, Section 5 summarizes the research and concludes the paper.

2 Related work

We outline the related work in this section. First, we give a short overview of privacy-preserving process mining in Section 2.1 and outline related work in the area of event log abstraction in Section 2.2.

2.1 Privacy-aware process mining

Differential privacy [6] has been widely studied in the process mining literature [11,14,28]. One big advantage of differential privacy is that it is immune to post-processing. This means that process models generated from anonymized event logs are also secure against attacks like the one introduced by Kirchmann et al. [19]. The literature knows two types of mechanisms for control-flow anonymization: They anonymize either the *directly-follows distribution* of an event log or the *trace-variant distribution*. Most anonymization techniques are based on noise insertion. While many techniques exist for anonymization, the impact of pre-processing was not studied in the literature. The only exception is work by Elkoumy et al. [8] that implemented an approach of privacy amplification. The idea is that through sub-sampling of traces a higher privacy guarantee can be given with less noise than if a differential privacy mechanism would be applied directly. However, this approach does not use process specific knowledge. Through our work, we present the first approach that utilizes event data partitioning in the pre-processing to achieve better utility for process discovery.

A direction differing from the work on differential privacy is focused on group-based privacy guarantees. The main line of research here, focuses on generalization [17], substitution [26], or the merging of events [13]. While we study our approach in combination with one differential privacy mechanism, it could also be applied in combination with one of these group-based privacy approaches.

Furthermore, non-anonymization approaches for privacy have been explored. These are usually focused on deriving process models in settings where the event logs are distributed over multiple data owners. These approaches use techniques such as multi-party computation [9] and trusted execution environments [15]. In contrast to these approaches, our technique can be applied to event logs stored by a single data owner.

2.2 Event log abstraction

Event log abstraction involves identifying sub-processes and their instances from low-level events [22]. Various studies have proposed approaches for this, often addressing the identification of sub-processes and their instances separately. For identifying sub-processes, Baier et al. [2] assumed that low-level activities with similar names originate from the same sub-process. Günther et al. [16] assumed that low-level activities appearing close together in an event log likely originate from the same higher-level activity. Lu et al. [23] proposed techniques that use domain knowledge and clustering.

Identifying sub-process instances has also been studied. A straightforward approach is to classify consecutive events with the same sub-process as originating from the same activity instance [16,7]. Other approaches include assuming that a sub-process can only be executed once [23] or setting a time limit for the execution of a single instance to distinguish between instances [2].

In conclusion, the privacy-aware process mining literature did not focus on how pre-processing can be helpful for better utility. However, it was shown in the literature that pre-processing such as event log abstraction can lead to better process discovery utility. In this work, we address this gap in the literature by studying how event log abstraction impacts privacy-aware process discovery.

3 Anonymization utilizing Event Data Partitioning

In Figure 2 we give an overview of the anonymization pipeline that allows us to combine event data partitioning and anonymization. In the remainder of this section, we elaborate on the different components of the pipeline.

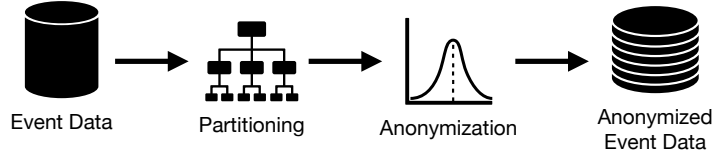


Fig. 2. Anonymization pipeline.

Event Data. Let us define an event as the execution of an activity $\alpha \in \mathcal{A}$, with \mathcal{A} being the universe of all activities. Let a trace $\sigma \in \mathcal{A}^*$ be a sequence of events. Each trace represents the execution of a process. The sequence of a trace indicates in which order the events have been executed. We define an event log $L \in \mathcal{L}$ as a multi-set of traces.

Partitioning. We turn to event data partitioning. First, we need to define event abstraction as a function $\psi : \mathcal{A} \mapsto \mathcal{A}$. Let \mathcal{A} always contain \top , which is the root activity, the highest-level activity within event abstraction. ψ allows us to abstract low-level activities, such as the reading of a sensor, into higher-level activities. Note that multiple low-level activities can be mapped onto the same higher-level activity. We assume that an activity a can not be abstracted to itself $\psi(a) \neq a$, with the exception of $\psi(\top) = \top$. The hierarchy ψ can be either user-defined or derived through automated approaches such as clustering. We assume that when multiple low-level activities are abstracted into the same higher-level activity α , the higher-level activity α represents a sub-process composed of these lower-level activities. A hierarchy ψ has to be defined such there are no cycles, e.g. $\psi(a) = \alpha$, $\psi(\alpha) = a$, and every activity can eventually be abstracted to the root activity \top , e.g. $\psi(a) = \alpha$, $\psi(\alpha) = \top$.

Our running example shown in Figure 3, illustrates such a process with such an event abstraction hierarchy. Here, the activities $\{a, b, c\}$ can be abstracted

to the activity A . Now we could apply the function ψ to every event within a trace σ to create an abstracted trace σ' . In such a scenario, all information that can be derived from the low-level activities would be lost. Therefore, it would be better to *decompose* a log L into multiple sub-logs: $\{L_\psi, L_{\alpha_1}, L_{\alpha_2}\}$ where L_ψ is the event log consisting of the higher-level activities while L_{α_1} contains all the activities belonging to sub-process α_1 . These composed logs can be generated by the *event data partitioning* function $\psi_L : \mathcal{L} \mapsto 2^{\mathcal{L}}$.

In algorithm 1, we outline how ψ_L works. First, we initialize the result set of sub-logs S and the abstracted event log L_ψ (see line 1-2). Next, we iterate over every trace σ in our event log L (see line 3). Next, we create the set C that is used to store traces of the sub-processes (see line 4). Now, we check for every activity of trace σ if it can be abstracted using ψ (see line 5-6). If an activity $\sigma(i)$ can be abstracted, we abstract it within the trace σ (see line 7 & line 10) and save the sub-process information belonging to the sub-process α' in a corresponding trace $\sigma'_{\alpha'}$ (see line 8-9). Such a procedure, leads to many occurrences of the sub-process activity α' in the trace σ . We apply a *compression function* that removes every occurrence except the first and last of activity α' from the trace σ , signaling the start and end of the sub-process (see line 13), while the actual information is encoded in the sub-process trace $\sigma'_{\alpha'}$. We add the trace ψ to the abstracted event log L (see line 14) and add the sub-process traces stored in the set C to the set of event logs S (see line 15).

Algorithm 1: Event Log Partitioning function ψ_L

Input: An event log L
Output: Set of event logs $S \in 2^{\mathcal{L}}$

```

1  $S \leftarrow \emptyset$ ; // Defines the empty set of new event logs
2  $L_\psi \leftarrow \emptyset$ ;
3 for  $\sigma \in L$  do
4    $C \leftarrow \emptyset$ ;
5   for  $i \in [1, \dots, |\sigma|]$  do
6     if  $\psi(\sigma(i)) \neq \emptyset$  then
7        $\alpha' \leftarrow \psi(\sigma(i))$ ;
8        $\sigma'_{\alpha'} \leftarrow C(\alpha')$  if  $\sigma_{\alpha'} \in C$ ; otherwise  $\sigma'_{\alpha'} \leftarrow \emptyset$ ;
9        $C \leftarrow (C \setminus \sigma_{\alpha'}) \cup \sigma'_{\alpha'} \circ \sigma(i)$ ;
10       $\sigma(i) \leftarrow \alpha'$ ;
11    end
12  end
13   $\text{compress}(\sigma)$ ;
14   $L_\psi \leftarrow L_\psi \cup \sigma$ ;
15   $S.\text{update}(C)$ ;
16 end
17  $S \leftarrow S \cup L_\psi$ ;
18 return  $S$ 
```

Anonymization. Next, we anonymize the partitioned event data. We apply the notion of differential privacy. We define $\gamma : \mathcal{L} \mapsto \mathcal{L}$ as an anonymization function that transforms an event log into an anonymized event log. Further, let

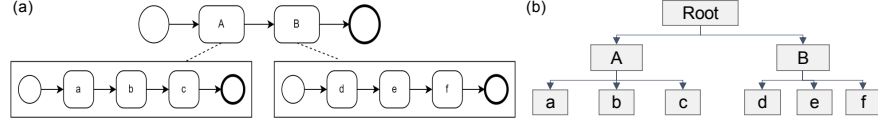


Fig. 3. Running example for process with sub-processes.

us define $img(\gamma) \subseteq \mathcal{L}$ as the *image* of γ , i.e., the set of all event logs that may be returned by γ . Finally, we define two event logs $L_1, L_2 \in \mathcal{L}$ to be *neighboring*, if they differ by exactly the data of one individual. In our setting, this corresponds to one trace, i.e., $|L_1 \setminus L_2| + |L_2 \setminus L_1| = 1$. Based on [6], we use the following definition for differential privacy:

Definition 1 (Differential Privacy). *Given an anonymization function γ and privacy parameter $\epsilon \in \mathbb{R}$, function γ provides ϵ -differential privacy, if for all neighboring pairs of event logs $L_1, L_2 \in \mathcal{L}$ and all subsets $\rho_\gamma \subseteq img(\alpha)$, it holds that:*

$$Pr[\gamma(L_1) \in \rho_\gamma] \leq e^\epsilon \times Pr[\gamma(L_2) \in \rho_\gamma]$$

where the probability is taken over the randomness introduced by the anonymization function γ .

We can apply any implementation of γ that gives the differential privacy guarantee. An important rule of differential privacy is parallel composition. This rule states that different datasets that have no intersection can be anonymized independently with ϵ -differential privacy their combination also is ϵ -differential private. Therefore, we can anonymize the different event logs created by the event log partitioning independently.

Anonymized Event Data. After anonymization, we end up with multiple anonymized event logs, i.e., L'_ψ , L'_A and L'_B . The anonymized log L'_ψ shows the process consisting of activities of higher level. The other logs contain the data of the anonymized sub-processes. For each of these logs, we apply a discovery algorithm and evaluate the quality of the discovered process model to measure the *utility*. Note that the other way around is also possible, namely, we first anonymize the event log and then partition the log into multiple sub logs. In the evaluation, we investigate the difference in the resulting utilities (RQ. 2).

4 Evaluation

In this section, we evaluate our proposed pipeline. Specifically, we examine the results for both directly-follows-based and trace-variant-based anonymization techniques. We focus specifically on the following two aspects:

- RQ1** How does performing partitioning before anonymization affect the utility of anonymized logs for process discovery?
- RQ2** How does the order of partitioning and anonymization impact process discovery utility?

In Section 4.1, we outline our experimental design and the datasets we employ. In Section 4.2, we present the results of our experiments. Finally, in Section 4.3, we provide a qualitative discussion on the trade-offs of combining partitioning with anonymization.

4.1 Experimental Setup

We provide an overview of our experimental setup in Figure 4. The first question (RQ1) aims to investigate the effect of our proposed pipeline (i.e., first partitioning and then anonymization, highlighted in orange in Figure 4) on utility (i.e., the quality of the process models discovered using anonymized logs). One may argue it is unfair to compare the utility of a set of logs to one log. Therefore, in the second question (RQ2), we examine whether any observed improvement is solely due to the introduction of partitioning rather than its specific placement in the pipeline. To investigate this, we compare the resulting utility of partitioning before (highlighted in orange) and after (highlighted in black) anonymization.

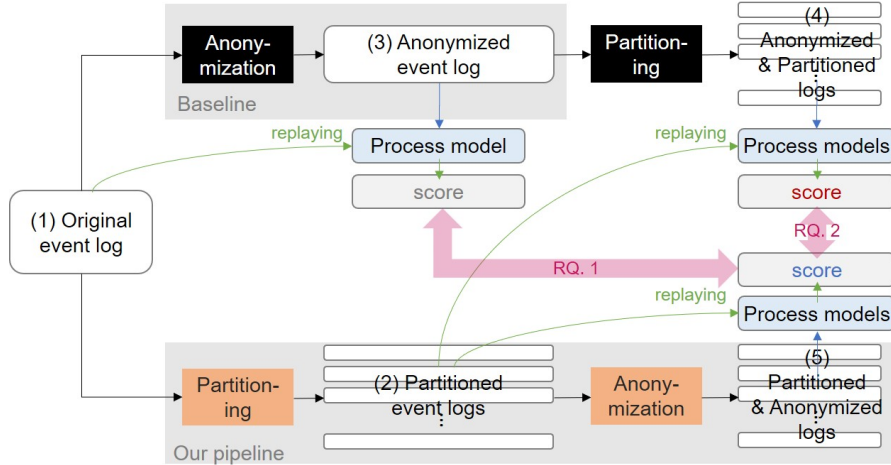


Fig. 4. Experimental setup overview.

Datasets. For our evaluation, we used the BPIC2012 [4], BPIC2015 [1]⁵, and BPIC2017 [5] event logs. All three event logs stem from real-world business processes. Both BPIC2012 and BPIC2017 describe a loan application process. BPIC2015 describes a permit application process used by Dutch municipalities. All of these logs are known to consist of events from multiple sub-processes and have therefore been widely used in various event abstraction studies [21,29,23]. In Table 1, we provide descriptive statistics of these datasets.

⁵ We used filtered log described in [1]. Although five datasets are available, we only used the first one, BPIC_1f.

Table 1. Statistical description of the datasets.

Data	#case	#events	#acts	avg. e/c	#trace variants
BPIC2012	13,087	262,200	24	20	4,366
BPIC2015	902	21,656	70	24	295
BPIC2017	31,509	1,202,267	26	38	15,930

Applied techniques for event data partitioning. During event data partitioning, we utilized three techniques for the higher-level activity mapping function (ψ in Section 3). The first technique developed by Lu et al. [23] integrates domain knowledge into event abstraction (FHM, FlexHMiner). Second, we use a technique by Günther et al. [16] based on activity clustering (AC). Finally, we employ a random clustering that allows us to study if we also experience benefits with less sophisticated abstraction techniques. We assumed that a high-level activity could only be executed once and distinguished between the start and end of higher-level activities, following the approach in [23].

Applied techniques for anonymization. For anonymization, we used two anonymization techniques: one based on the Laplace mechanism (DF-Laplace) [24] (since it does not require an additional parameter, as does the technique based on the exponential mechanism such as SaPa [14]); the other one based on SaCoFa [14] for the trace-variant query, since it provides the highest utility of all techniques that add additional behavior to anonymized logs. We used public implementations of both techniques [18]. For differential privacy, we used the settings of $\varepsilon = \{0.01, 0.1, 1.0\}$ to test different orders of magnitudes for privacy guarantee. Lower values of ε mean stronger privacy protection. In our result figures, we plot the results of these settings within one bin, due to space limitations.

SaCoFa anonymizes trace-variants by constructing a prefix tree that accounts for the harmfulness of the prefix. The upper bound on the trace-variant length (l) was set to 50 across all settings. However, the privacy guarantee (ε) and the pruning parameter (p) were evaluated at three different combinations: (0.01, 300), (0.1, 200), and (1.0, 100) for BPIC2012 and BPIC2017, and (0.01, 100), (0.1, 100), and (1.0, 50) for BPIC2015. To account for non-determinism of noise injection, we repeated the anonymization process 10 times for each setting.

Applied techniques for process discovery. For process discovery, we used Inductive Miner [20] and Heuristic Miner [33], setting the noise threshold to 0.2 for both techniques. Both techniques are widely used process discovery algorithms and we used their implementation in PM4Py [3].

Evaluation measures. We evaluated the utility of the process models using four metrics. For fitness, we utilized token-based replay fitness [31]. For precision, we employed ETC Precision [25]. We also calculated the F1-score, which is the harmonic mean of fitness and precision. Finally, for generalization, we adapted the technique from [1] with $k = 3$.

When evaluating the multi-level process model, we followed the approach in [23]. This involves calculating the performance for each sub-process in the multi-level process model and then averaging them.

4.2 Results

RQ1: Effect of Partitioning on DF-Laplace. Figure 5 compares the utility of performing only anonymization with that of performing anonymization after partitioning, categorized by dataset, applied technique in event partitioning, and process discovery technique. The fitness was higher when only anonymization was applied, for almost all datasets. On the other hand, precision exhibited significant benefits from event partitioning: it was below 0.2 when only anonymization was performed, but increased to as high as 0.6 after partitioning at most. Overall, the increase in precision contributed to higher F1-score values when anonymization was preceded by event partitioning. Generalization showed good utility in most cases, with values exceeding 0.8 regardless of whether partitioning was applied, and differences remained within 0.2. These results indicate that incorporating partitioning, particularly with the DF-Laplace anonymization technique, can enhance precision while maintaining good levels of fitness and generalization.

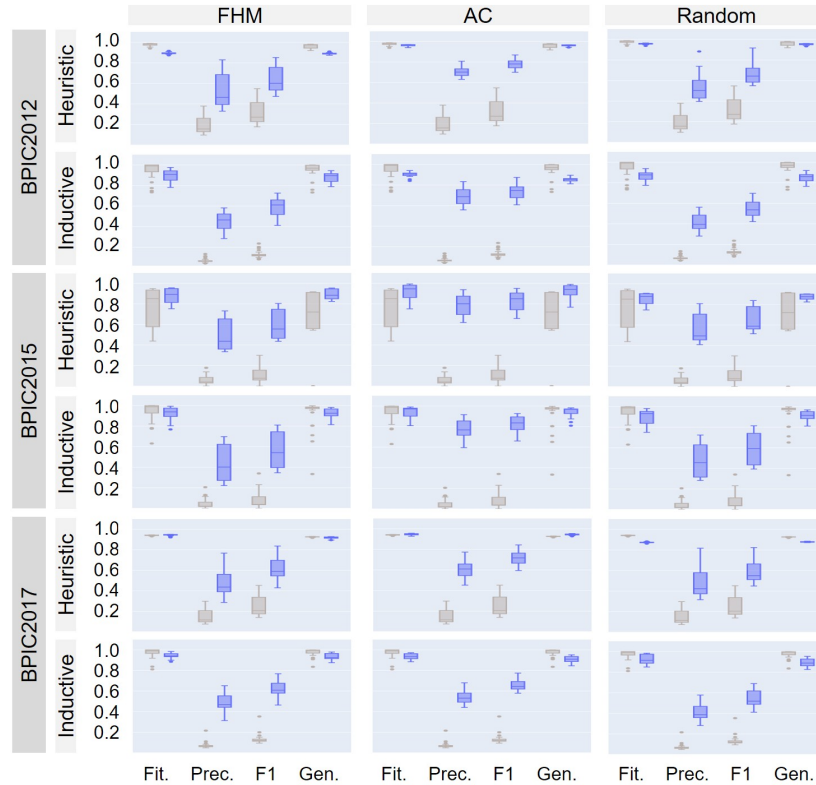


Fig. 5. Utility comparison with DF-Laplace: Anonymization only (gray) vs. Partitioning and Anonymization (blue)

RQ1: Effect of Partitioning on SaCoFa. When SaCoFa was used for anonymization, the results differed from those observed with DF-Laplace. While partitioning tended to improve the utility of anonymized logs when DF-Laplace was used, it appeared to have little to no effect on utility when SaCoFa was applied. As shown in Figure 6, there was generally little difference across all metrics between applying anonymization alone and applying partitioning before anonymization. Specifically, when DF-Laplace was used, precision showed a significant improvement after partitioning, but with SaCoFa, the difference in precision between logs with and without partitioning was minimal.

This outcome can be attributed to the fact that DF-Laplace significantly reduces precision when applied to the full log, whereas SaCoFa maintains relatively stable precision levels, and in some cases even improves them.

For the BPIC2017 dataset, however, precision showed an improvement compared to other datasets. This can be explained by the relatively low precision of the BPIC2017 log after anonymization, allowing partitioning to have a more noticeable effect.

For the BPIC2015 dataset, precision itself did not improve significantly, but partitioning effectively reduced its variation. This suggests that while anonymization alone results in good precision on average, there are cases where its quality is not consistently guaranteed. In such instances, partitioning can help stabilize the quality of anonymization and ensure more reliable results.

RQ2: Effect of Partitioning Order on DF-Laplace. Figure 7 compares the utility of performing partitioning after anonymization and the utility of performing partitioning before anonymization, categorized by dataset, applied abstraction technique in event partitioning, and process discovery technique. Overall, across all datasets, performing partitioning before anonymization yielded better performance in all metrics compared to partitioning after anonymization. However, fitness and generalization generally exceeded 0.8 in both cases, resulting in relatively minor differences, while precision showed a more pronounced disparity. For example, when partitioning was performed using random clustering on BPIC2015 and the process model was discovered using the heuristic algorithm, the precision improved significantly—from approximately 0.3 when partitioning was performed after anonymization to 0.8 when partitioning was performed beforehand. In conclusion, these results demonstrate that when using DF-Laplace for anonymization, partitioning before anonymization helps to preserve more information in the data, thereby enhancing process discovery performance.

When partitioning was performed using FHM on BPIC2015, the results for precision differed somewhat from other observations. While precision generally improved significantly when partitioning was performed first, in this case it decreased. Considering that the precision was high when FHM was performed on BPIC2015 without anonymization (0.98 with the heuristic algorithm and 0.89 with the inductive algorithm), this result can be interpreted as the impact of data distortion from anonymization being more pronounced in the sub-logs.

RQ2: Effect of Partitioning Order on SaCoFa. For BPIC2015 and BPIC2017, partitioning first showed slightly better performance, but overall the differences

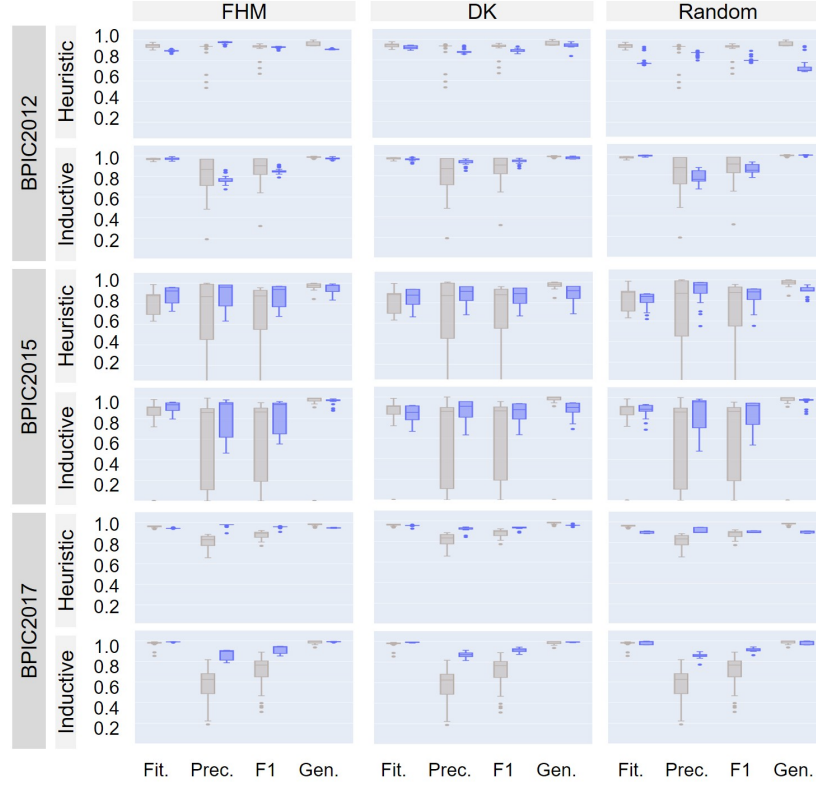


Fig. 6. Utility comparison with SaCoFa: Anonymization only (gray) vs. Partitioning and Anonymization (blue).

were minimal. In BPIC2012, partitioning later performed slightly better. Looking back at the earlier results with DF-Laplace, where applying partitioning last significantly reduced precision—leading to much better performance when partitioning was done beforehand—this outcome can be interpreted differently for SaCoFa. Since SaCoFa does not significantly degrade precision even when anonymization is applied without partitioning, this likely explains the observed results.

However, when applying FHM to BPIC2015, partitioning first resulted in significantly better utility compared to partitioning later. This is an unusual result, considering that for the same dataset, when using Günther’s technique or Random clustering, there was little difference between partitioning first and partitioning later. Comparing this with the results of using SaCoFa for anonymization alone (Figure 6), we see that when FHM was applied to anonymized data, precision dropped significantly. In contrast, when other techniques were applied, the precision remained relatively stable. This suggests that the effectiveness of an abstraction technique used for partitioning can vary.

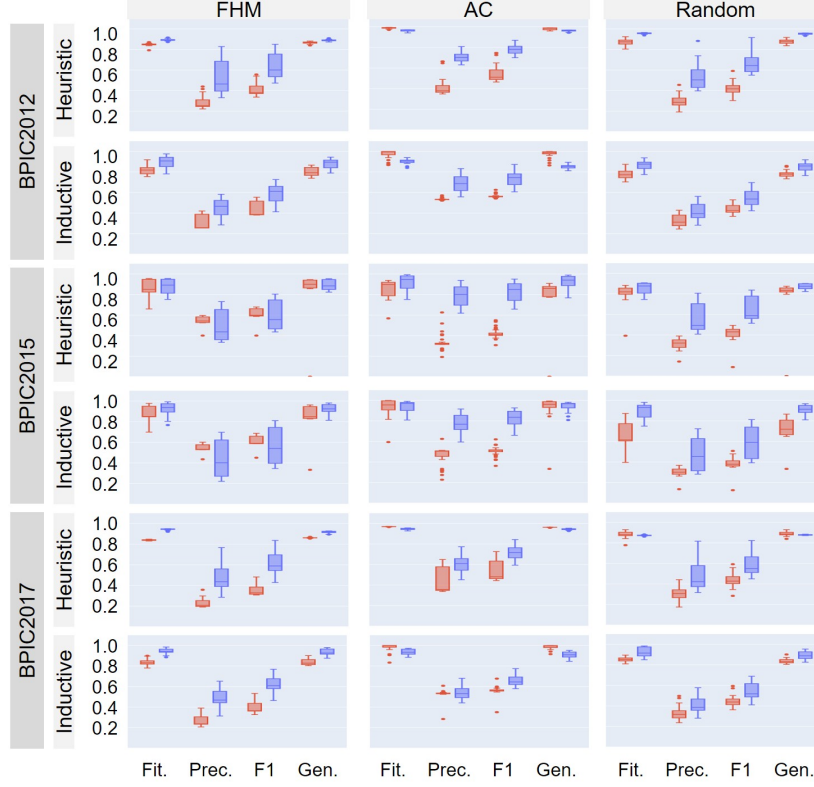


Fig. 7. Utility comparison with DF-Laplace: Anonymization and Partitioning (red) vs. Partitioning and Anonymization (blue).

4.3 Trade-Off Discussion

In this section, we discuss trade-offs [30] that come from integrating event data partitioning into the anonymization process.

Loss of Information due to Partitioning. Event log partitioning can remove information that is useful for the analysis of the business process. Therefore, event partitioning can lead to an additional information loss that is difficult to quantify. This trade-off has to be considered when weighing the utility gains from lower noise addition.

Explosion of Pipeline Options. Choosing the right anonymization technique is challenging as there is no standardized guide. This creates a risk that users may apply a suboptimal anonymization technique for their analytical needs. In this paper, we introduce a pipeline that further increases the complexity of this decision. Instead of selecting only an anonymization technique, users must also choose an abstraction technique for event data partitioning. As a result, one drawback of our pipeline is that identifying the optimal configuration for maximizing utility in each scenario requires substantial domain expertise.

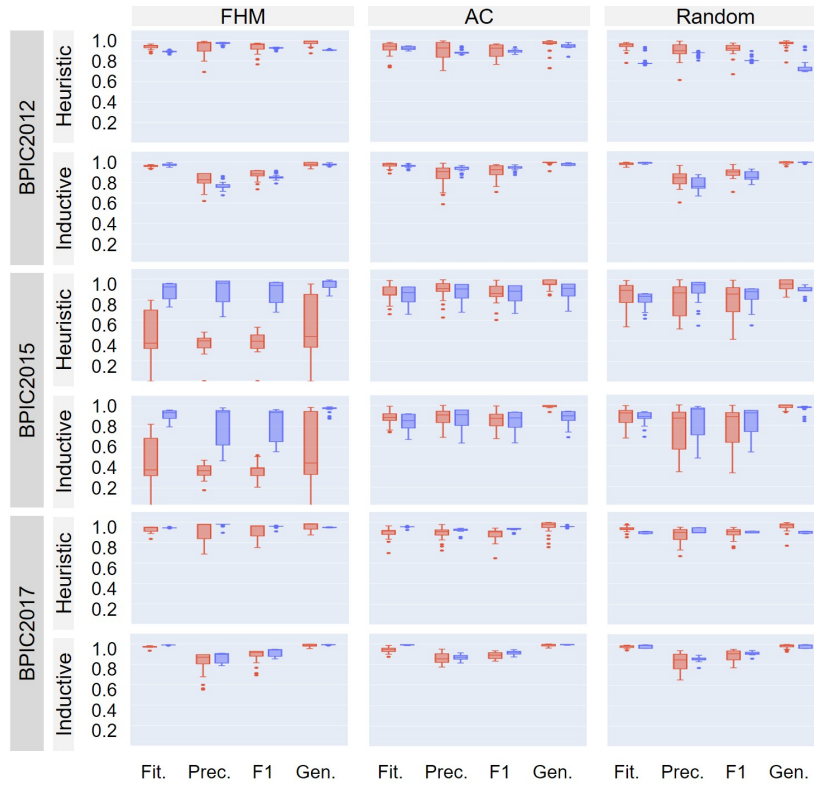


Fig.8. Utility comparison with SaCoFa: Anonymization and Partitioning (red) vs. Partitioning and Anonymization (blue).

Less Uncertainty of the Sub-processes. When anonymizing the whole process, a lot of behavior can be added as noise. However, when sub-processes are anonymized independently, the new behavior that can be added is significantly restricted. This may lead an adversary to learn more information about a sub-process than would be possible using anonymization without event partitioning. It is important to remember that differential privacy guarantees that the same privacy protection is given to individuals either way.

5 Conclusion

In this paper, we address the problem of privacy-aware process discovery in a novel way. We propose a pipeline that builds on event data partitioning before anonymization. To evaluate our pipeline, we investigated the impact of event partitioning on anonymization for automatic process discovery. Our findings demonstrated that our pipeline provides outperformance when directly-follows-based anonymization techniques are applied. We believe this findings suggest

that combining anonymization with event log pre-processing has the opportunity to unlock higher utility in privacy-aware process mining. For future work, we intend to explore how other pre-processing techniques can contribute to better utility. Ideally, we aim to develop a framework that helps determine the most effective combination of pre-processing and anonymization techniques based on event log characteristics.

References

1. Augusto, A., Conforti, R., Dumas, M., Rosa, M.L., Maggi, F.M., Marrella, A., Mecella, M., Soo, A.: Automated discovery of process models from event logs: Review and benchmark. *IEEE Trans. Knowl. Data Eng.* **31**(4), 686–705 (2019)
2. Baier, T., Mendling, J., Weske, M.: Bridging abstraction layers in process mining. *Information Systems* **46**, 123–139 (2014)
3. Berti, A., Van Zelst, S.J., van der Aalst, W.: Process mining for python (pm4py): bridging the gap between process-and data science. *arXiv preprint arXiv:1905.06169* (2019)
4. van Dongen, B.: Bpi challenge 2012 (2012), https://data.4tu.nl/articles/dataset/BPI_Challenge_2012/12689204/1
5. van Dongen, B.: Bpi challenge 2017 (2017), https://data.4tu.nl/articles/dataset/BPI_Challenge_2017/12696884/1
6. Dwork, C.: Differential privacy. In: *International colloquium on automata, languages, and programming*. pp. 1–12. Springer (2006)
7. van Eck, M.L., Sidorova, N., van der Aalst, W.M.P.: Enabling process mining on sensor data from smart products. In: *RCIS*. pp. 1–12. IEEE (2016)
8. Elkoumy, G., Dumas, M.: Libra: High-utility anonymization of event logs for process mining via subsampling. In: Burattin, A., Polyvyanyy, A., Weber, B. (eds.) *ICPM*. pp. 144–151. IEEE (2022)
9. Elkoumy, G., Fahrenkrog-Petersen, S.A., Dumas, M., Laud, P., Pankova, A., Weidlich, M.: Secure multi-party computation for inter-organizational process mining. In: *BPMDS*. Springer (2020)
10. Elkoumy, G., Fahrenkrog-Petersen, S.A., Sani, M.F., Koschmider, A., Mannhardt, F., Von Voigt, S.N., Rafiei, M., Waldthausen, L.V.: Privacy and confidentiality in process mining: Threats and research challenges. *ACM TMIS* (2021)
11. Elkoumy, G., Pankova, A., Dumas, M.: Differentially private release of event logs for process mining. *Inf. Syst.* **115**, 102161 (2023)
12. Fahrenkrog-Petersen, S.A., van der Aa, H., Weidlich, M.: PRIPEL: privacy-preserving event log publishing including contextual information. In: *BPM. Lecture Notes in Computer Science*, vol. 12168, pp. 111–128. Springer (2020)
13. Fahrenkrog-Petersen, S.A., van der Aa, H., Weidlich, M.: Optimal event log sanitization for privacy-preserving process mining. *DKE* **145**, 102175 (2023)
14. Fahrenkrog-Petersen, S.A., Kabierski, M., van der Aa, H., Weidlich, M.: Semantics-aware mechanisms for control-flow anonymization in process mining. *Information Systems* **114**, 102169 (2023)
15. Goretta, V., Basile, D., Barbaro, L., Ciccio, C.D.: Trusted execution environment for decentralized process mining. In: Guizzardi, G., Santoro, F.M., Mouratidis, H., Soffer, P. (eds.) *CAiSE* (2024)
16. Günther, C.W., Rozinat, A., van der Aalst, W.M.P.: Activity mining by global trace segmentation. In: *BPM Workshops* (2009)

17. Hildebrant, R., Fahrenkrog-Petersen, S.A., Weidlich, M., Ren, S.: PMDG: privacy for multi-perspective process mining through data generalization. In: CAiSE. Springer (2023)
18. Kirchmann, H., Fahrenkrog-Petersen, S.A., Kabierski, M., van der Aa, H., Weidlich, M.: Privacy-preserving process mining with pm4py (extended abstract). In: ICPM Demo. CEUR-WS.org (2022)
19. Kirchmann, H., Fahrenkrog-Petersen, S.A., Mannhardt, F., Weidlich, M.: Control-flow reconstruction attacks on business process models. CoRR **abs/2409.10986** (2024), <https://doi.org/10.48550/arXiv.2409.10986>
20. Leemans, S.J., Fahland, D., Van Der Aalst, W.M.: Discovering block-structured process models from event logs containing infrequent behaviour. In: BPM Workshops. Springer (2014)
21. Li, C.Y., van Zelst, S.J., van der Aalst, W.M.: Event abstraction for partial order patterns. In: BPM. pp. 38–54. Springer (2023)
22. Lim, J., Song, M.: Navigating event abstraction in process mining: A comprehensive analysis of sub-problems, data, and process characteristic considerations. In: BPM Workshops. pp. 174–185. Springer (2024)
23. Lu, X., Gal, A., Reijers, H.A.: Discovering hierarchical processes using flexible activity trees for event abstraction. In: ICPM. pp. 145–152. IEEE (2020)
24. Mannhardt, F., Koschmider, A., Baracaldo, N., Weidlich, M., Michael, J.: Privacy-preserving process mining: Differential privacy for event logs. BISE **61**, 595–614 (2019)
25. Munoz-Gama, J., Carmona, J.: A fresh look at precision in process conformance. In: BPM. pp. 211–226. Springer (2010)
26. Rafiei, M., van der Aalst, W.M.P.: Group-based privacy preservation techniques for process mining. DKE **134**, 101908 (2021)
27. Rafiei, M., van der Aalst, W.M.: Towards quantifying privacy in process mining. In: ICPM Workshops. pp. 385–397. Springer (2021)
28. Rafiei, M., Wangelik, F., Pourbafrani, M., van der Aalst, W.M.P.: Travag: Differentially private trace variant generation using gans. In: RCIS. Lecture Notes in Business Information Processing, vol. 476, pp. 415–431. Springer (2023)
29. Rebmann, A., Weidlich, M., van der Aa, H.: GECCO: constraint-driven abstraction of low-level event logs. In: ICDE. pp. 150–163. IEEE (2022)
30. Robillard, M.P., Arya, D.M., Ernst, N.A., Guo, J.L., Lamothe, M., Nassif, M., Novielli, N., Serebrenik, A., Steinmacher, I., Stol, K.J.: Communicating study design trade-offs in software engineering. ACM TOSEM (2024)
31. Van Der Aalst, W., van der Aalst, W.: Data science in action. Springer (2016)
32. von Voigt, S.N., Fahrenkrog-Petersen, S.A., Janssen, D., Koschmider, A., Tschorsch, F., Mannhardt, F., Landsiedel, O., Weidlich, M.: Quantifying the re-identification risk of event logs for process mining. In: CAiSE. Springer (2020)
33. Weijters, A., Ribeiro, J.: Flexible heuristics miner (fhm). In: CIDM. pp. 310–317 (2011)