

Preserving Data Secrecy in Decentralized Inter-organizational Process Mining

Valerio Goretti^a, Davide Basile^a, Luca Barbaro^a, Claudio Di Ciccio^b

^a*Sapienza University of Rome, Via Regina Elena 295, 00161, Rome, Italy*

^b*Utrecht University, Princetonplein 5, 3584 CC, Utrecht, The Netherlands*

Abstract

Inter-organizational business processes involve multiple independent organizations collaborating to achieve mutual interests. Process mining techniques have the potential to allow these organizations to enhance operational efficiency, improve performance, and deepen the understanding of their business based on the recorded process event data. However, inter-organizational process mining faces substantial challenges, including topical secrecy concerns: The involved organizations may not be willing to expose their own data to run mining algorithms jointly with their counterparts or third parties. In this paper, we introduce CONFINE, a novel approach that unlocks process mining on multiple actors' process event data while safeguarding the secrecy and integrity of the original records in an inter-organizational business setting. To ensure that the phases of the presented interaction protocol are secure and that the processed information is hidden from involved and external actors alike, our approach resorts to a decentralized architecture comprised of trusted applications running in Trusted Execution Environments (TEEs). We show the feasibility of our solution by showcasing its application to a healthcare scenario and evaluating our implementation in terms of memory usage and scalability on real-world event logs.

Keywords: Collaborative information

1. Introduction

In today’s business landscape, organizations constantly seek ways to enhance operational efficiency, increase performance, and gain valuable insights to improve their processes. Process mining offers techniques to discover, monitor, and improve business processes by extracting knowledge from chronological records known as *event logs* [1]. Process-aware information systems record events referring to activities and interactions within a business process. The vast majority of process mining contributions consider *intra-organizational* settings, in which business processes are executed inside individual organizations. However, organizations increasingly recognize the value of collaboration and synergy in achieving operational excellence. *Inter-organizational* business processes involve several independent organizations cooperating to achieve a shared objective [2]. Despite the advantages of transparency, performance optimization, and benchmarking that companies can gain, inter-organizational process mining raises challenges that hinder its application. The major issue concerns confidentiality. Companies are reluctant to share private information required to execute process mining algorithms with their partners [3]. Indeed, letting sensitive operational data traverse organizational boundaries introduces concerns about data privacy, security, and compliance with internal regulations [4]. *Trusted Execution Environments* (TEEs) [5] can serve as fundamental enablers to balance the need for insights with the need to protect sensitive information in inter-organizational settings. TEEs offer secure contexts that guarantee code integrity and data confidentiality before, during, and after its utilization.

In this paper, we propose CONFINE, a novel approach and tool aimed at enhanc-

ing collaborative information system architectures with secrecy-preserving process mining capabilities in a decentralized fashion. It resorts to *trusted applications* running in TEEs to preserve the secrecy and integrity of shared data. To pursue this aim, we design a decentralized architecture for a four-staged protocol: (i) The initial exchange of preliminary metadata, (ii) the attestation of the miner entity, (iii) the secure transmission and privacy-preserving merge of encrypted information segments amid multiple parties, (iv) the isolated and verifiable computation of process discovery algorithms on joined data. We evaluate our proof-of-concept implementation against synthetic and real-world-based data with a convergence test followed by experiments to assess the scalability of our approach.

The remainder of this paper is as follows. [Sect. 2](#) provides an overview of related work. In [Sect. 3](#), we introduce a motivating use-case scenario in healthcare. We present the CONFINE approach in [Sect. 5](#). We describe the implementation of our approach in [Sect. 6](#). In [Sect. 7](#), we report on the efficacy and efficiency tests for our solution. Finally, we conclude our work and outline future research directions in [Sect. 8](#).

2. Background and Related Work

2.1. Background

2.1.1. Inter-organizational Process Mining

2.1.2. Trusted Execution Environments

2.2. Related Work

Despite the relative recency of this research branch across process mining and collaborative information systems, scientific literature already includes noticeable contributions to inter-organizational process mining. The work of Müller et al. [\[6\]](#)

47 focuses on data privacy and security within third-party systems that mine data gener-
48 ated from external providers on demand. To safeguard the integrity of data earmarked
49 for mining purposes, their research introduces a conceptual architecture that entails
50 the execution of process mining algorithms within a cloud service environment,
51 fortified with Trusted Execution Environments. Drawing inspiration from this foun-
52 dational contribution, our research work seeks to design a decentralized approach
53 characterized by organizational autonomy in the execution of process mining algo-
54 rithms, devoid of synchronization mechanisms involvement taking place between the
55 involved parties. A notable departure from the framework of Müller et al. lies in the
56 fact that here each participating organization retains the discretion to choose when
57 and how mining operations are conducted. Moreover, we bypass the idea of fixed
58 roles, engineering a peer-to-peer scenario in which organizations can simultaneously
59 be data provisioners or miners. Elkoumy et al. [7, 8] present Shareprom. Like our
60 work, their solution offers a means for independent entities to execute process mining
61 algorithms in inter-organizational settings while safeguarding their proprietary input
62 data from exposure to external parties operating within the same context. Share-
63 prom’s functionality, though, is confined to the execution of operations involving
64 event log abstractions [9] represented as directed acyclic graphs, which the parties
65 employ as intermediate pre-elaboration to be fed into secure multiparty computation
66 (SMPC) [10]. As the authors remark, relying on this specific graph representation
67 imposes constraints that may prove limiting in various process mining scenarios. In
68 contrast, our approach allows for the secure, ciphered transmission of event logs to
69 process mining nodes as a whole. Moreover, SMPC-based solutions require compu-
70 tationally intensive operations and synchronous cooperation among multiple parties,
71 which make these protocols challenging to manage as the number of participants

72 scales up [11]. In our research work, individual computing nodes run the calculations,
73 thus not requiring synchronization with other machines once the input data is loaded.

74 We are confronted with the imperative task of integrating event logs originat-
75 ing from different data sources and reconstructing consistent traces that describe
76 collaborative process executions. Consequently, we engage in an examination of
77 methodologies delineated within the literature, each of which offers insights into
78 the merging of event logs within inter-organizational settings. The work of Claes
79 et al. [12] holds particular significance for our research efforts. Their seminal study
80 introduces a two-step mechanism operating at the structured data level, contingent
81 upon the configuration and subsequent application of merging rules. Each such rule
82 indicates the relations between attributes of the traces and/or the activities that must
83 hold across distinct traces to be combined. In accordance with their principles, our
84 research incorporates a structured data-level merge based on case references and
85 timestamps as merging attributes. The research by Hernandez et al. [13] posits a
86 methodology functioning at the raw data level. Their approach represents traces
87 and activities as *bag-of-words* vectors, subject to cosine similarity measurements
88 to discern links and relationships between the traces earmarked for combination. An
89 appealing aspect of this approach lies in its capacity to generalize the challenge of
90 merging without necessitating a-priori knowledge of the underlying semantics inher-
91 ent to the logs under consideration. However, it entails computational overhead in
92 the treatment of data that can interfere with the overall effectiveness of our approach.

93 $T_{312} = \langle \text{PH, COPA, OD, DOR, PDL, SD, RD, AD, TP, PAFH, PIA, PT, VRT, TPB, RPB, DPH, PCD, DP} \rangle$

$T_{711} = \langle \text{PH, COPA, OD, DOR, PDL, SD, RD, AD, TP, DPH, PCD, DP} \rangle$

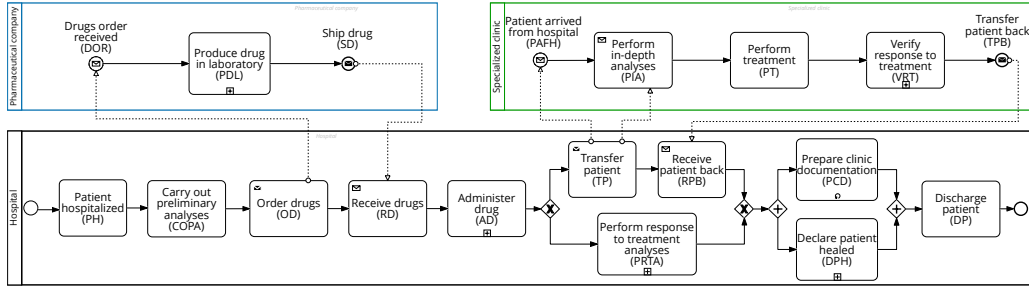


Figure 1: A BPMN collaboration diagram of a simplified healthcare scenario

Table 1: Events from cases 312 (Alice) and 711 (Bob) recorded by the hospital, the specialized clinic, and the pharmaceutical company

| Hospital | | | | | |
|----------|------------------|----------|------|------------------|----------|
| Case | Timestamp | Activity | Case | Timestamp | Activity |
| 312 | 2022-07-14T10:36 | PH | 312 | 2022-07-15T22:06 | TP |
| 312 | 2022-07-14T16:36 | COPA | 711 | 2022-07-16T00:55 | PRTA |
| 711 | 2022-07-14T17:21 | PH | 711 | 2022-07-16T00:55 | PCD |
| 312 | 2022-07-14T17:36 | OD | 711 | 2022-07-16T02:55 | DPH |
| 711 | 2022-07-14T23:21 | COPA | 711 | 2022-07-16T04:55 | DP |
| 711 | 2022-07-15T00:21 | OD | 312 | 2022-07-16T07:06 | RPB |
| 711 | 2022-07-15T18:55 | RD | 312 | 2022-07-16T09:06 | DPH |
| 312 | 2022-07-15T19:06 | RD | 312 | 2022-07-16T09:06 | PCD |
| 711 | 2022-07-15T20:55 | AD | 312 | 2022-07-16T11:06 | DP |
| 312 | 2022-07-15T21:06 | AD | 312 | 2022-07-16T11:06 | DP |

| Pharmaceutical company | | |
|------------------------|------------------|----------|
| HospitalCaseId | Timestamp | Activity |
| 312 | 2022-07-15T09:06 | DOR |
| 711 | 2022-07-15T09:30 | DOR |
| 312 | 2022-07-15T11:06 | PDL |
| 711 | 2022-07-15T11:30 | PDL |
| 312 | 2022-07-15T13:06 | SD |
| 711 | 2022-07-15T13:30 | SD |

| Specialized clinic | | |
|--------------------|------------------|----------|
| TreatmentID | Timestamp | Activity |
| 312 | 2022-07-16T00:06 | PAFH |
| 312 | 2022-07-16T01:06 | PIA |
| 312 | 2022-07-16T03:06 | PT |
| 312 | 2022-07-16T04:06 | VRT |
| 312 | 2022-07-16T05:06 | TPB |

3. Motivating Scenario

For our motivating scenario, we focus on a simplified hospitalization process for the treatment of rare diseases. The process model is depicted as a BPMN diagram in Fig. 1 and involves the cooperation of three parties: a hospital, a pharmaceutical company, and a specialized clinic. For the sake of simplicity, we describe the process through two cases, recorded by the information systems as in Table 2. Each patient in the hospital is associated with an id which would be the identifier of the case in the hospital log. Alice's journey (**case 312**) begins when she enters the hospital for the preliminary examinations (patient hospitalized, PH). The hospital then places an order for the drugs (OD) to the pharmaceutical company for treating Alice's specific condition. Afterwards, the pharmaceutical company acknowledges that the drugs

105 order is received (DOR), proceeds to produce the drugs in the laboratory (PDL),
 106 and ships the drugs (SD) back to the hospital. Upon receiving the medications, the
 107 hospital administers the drug (AD), and conducts an assessment to determine if
 108 Alice can be treated internally. If specialized care is required, the hospital transfers
 109 the patient (TP) to the specialized clinic. When the patient arrives from the hospital
 110 (PAFH), the specialized clinic performs in-depth analyses (PIA) and proceeds with
 111 the treatment (PT). Once the specialized clinic had completed the evaluations and
 112 verified the response to the treatment (VRT), it transfers the patient back (TPB). The
 113 hospital receives the patient back (RPB) and prepares the clinical documentation
 114 (PCD). If Alice has successfully recovered, the hospital declares the patient as healed
 115 (DPH). When Alice's treatment is complete, the hospital discharges the patient (DP).
 116 Bob (**case 711**) enters the hospital a few hours later than Alice. His hospitalization
 117 process is similar to Alice's. However, he does not need specialized care, and his
 118 case is only treated by the hospital. Therefore, the hospital performs the response
 119 to treatment analyses (PRTA) instead of transferring him to the specialized clinic.

120 Both the National Institute of Statistics of the country in which the three organiza-
 121 tions reside and the University that hosts the hospital wish to uncover information on
 122 this inter-organizational process for reporting and auditing purposes [14] via process
 123 analytics. The involved organizations share the urge for such an analysis and wish to
 124 be able to repeat the mining task also in-house. The hospital, the specialized clinic,
 125 and the pharmaceutical company have a partial view of the overall unfolding of the
 126 inter-organizational process as they record the events stemming from the parts of their
 127 pertinence. In Sect. 2.2, we show the cases 312 and 711 and the corresponding traces
 128 recorded by the hospital (i.e., T_{312} and T_{711}). The specialized clinic and the pharma-
 129 ceutical company have their internal ids. However, interactions between the hospital,

and the pharmaceutical company and the specialized clinic always have references to the hospital case id of the patients involved. In Table 1 you can see references to hospital ids in the first column. Since this link is present, it is not necessary that the case ids of pharmaceutical company and specialized clinic are synchronized with the hospital cases id. Those traces are projections of the two combined ones for the whole inter-organizational process: $T_{312} = \langle \text{PH, COPA, OD, DOR, PDL, SD, RD, AD, TP, PAFH, PIA, PT, VRT, TPB, RPB, DPH, PCD, DP} \rangle$ and $T_{711} = \langle \text{PH, COPA, OD, DOR, PDL, SD, RD, AD, TP, DPH, PCD, DP} \rangle$. Results stemming from the analysis of the local cases would not provide a full picture. Data should be merged. However, to preserve the privacy of the people involved and safeguard the confidentiality of the information, the involved parties cannot give open access to their traces to other organizations. The diverging interests (being able to conduct process mining on data from multiple sources without giving away the local event logs in-clear) motivate our research. In the following, we describe the design of our solution.

4. Preliminaries

Given a finite set of events \hat{E} and a total-order relation \leq subset of $\hat{E} \times \hat{E}$, we identify an event log as the totally ordered set (\hat{E}, \leq) . In the example, ... Let \widehat{IID} be a finite non-empty set of symbols such that $|\widehat{IID}| \leq |\hat{E}|$. We assume that every event be associated with a *case identifier* $iid \in \widehat{IID}$ via a total surjective function $iid : \hat{E} \rightarrow \widehat{IID}$ such that the restriction $<_{iid} = \leq \cap \{e \in \hat{E} : iid(e) = iid\}^2$ of total order \leq on all events mapped to the same iid is strict (i.e., if $e \leq e'$ with $e \neq e'$ and $iid(e) = iid(e')$ then $e' \not\leq e$). In the example, ... In other words, iid acts as an equivalence relation partitioning \hat{E} into $\{\hat{E}_{iid}\}_{iid \in \widehat{IID}} \subseteq 2^{\hat{E}}$ based on the iid to which the events $e \in \hat{E}_{iid}$ map, and imposing that events are linearly ordered by the restriction of \leq on every

Add example

Add example

\hat{E}_{iid} . Every pair $(\hat{E}_{iid}, <_{iid})$ thus represents a finite linearly totally ordered set (or
 loset for brevity) with $\hat{E}_{iid} \subseteq \hat{E}$ and $<_{iid} \subseteq \hat{E}_{iid} \times \hat{E}_{iid} \subseteq \leq \subseteq \hat{E} \times \hat{E}$. Let $(\hat{E}, <)$ be a
 loset and $(\hat{E}', <')$, $(\hat{E}'', <'')$ two (sub-)losets such that $\hat{E}' \cup \hat{E}'' \subseteq \hat{E}$ and $\hat{E}' \cap \hat{E}'' = \emptyset$,
 with $<'$ and $<''$ being the restrictions of $<$ on \hat{E}' and \hat{E}'' , respectively. We define the
 order-preserving union $\oplus: \hat{E}^3 \times \hat{E}^3 \rightarrow \hat{E}^3$ of losets as follows: $(\hat{E}', <') \oplus (\hat{E}'', <'') =$
 $(\hat{E}' \cup \hat{E}'', < \cap (\hat{E}' \cup \hat{E}'')^2)$. We can thus derive the notion of case C_{iid} given a $iid \in \widehat{IID}$
 as a loset of events mapping to the same iid and ordered by the linear restriction
 $<$ of \leq over the events in C_{iid} : $iid = (\hat{E}_{iid}, <)$ where $C_{iid} = \langle e_1, \dots, e_{|C_{iid}|} \rangle$ where
 $iid(e_i) = iid \in \widehat{IID}$ for every i s.t. $1 \leq i \leq |C_{iid}|$ and $e_i < e_j$ for every $i \leq j \leq |C_{iid}|$.¹
 Notice that the cardinality of \hat{C} and \widehat{IID} coincide. Events are also the domain of a
 function $p: \hat{E} \rightarrow \hat{\mathcal{P}}$ mapping events to log provisioners. In the example, . . . We shall
 denote with $C_{iid}^{\mathcal{P}}$ the loset consisting of every event $e \in C_{iid}$ such that $p(e) = \mathcal{P}$, with
 the restriction of the strict total order of C_{iid} on those events. In the example, . . .

Continue
here re-
vising the
definition
of order
preserving
union

Add exam-
ple

Add exam-
ple

Gotta move this one earlier when we introduce the example with the partitioned event log. (Section 4.2??). Clarify
 the difference between segmentation (given a segsize, i.e., a segment of a case-part in a sublog) and partitioning
 (of a log into case-parts of sublogs. Then, prove that the pipeline of partitioning and segmentation has its inverse
 in the union and merge for soundness.

5. Design

In this section, we present the high-level architecture of the CONFINE framework.
 We consider the main functionalities of each component, avoiding details on the
 employed technologies discussed in the next sections.

¹We employ the angular-bracket notation here for the sake of simplicity, although it is typically
 used for sequences. Unlike sequences, cases do not allow for the same event to occur more than once.

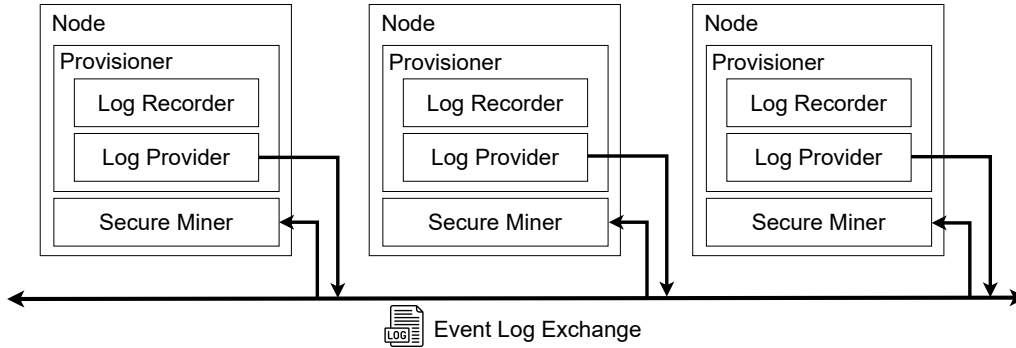


Figure 2: The CONFINE high-level architecture

172 **The CONFINE architecture at large.** Our ar-
 173 chitecture involves different information systems
 174 running on multiple machines. An organiza-
 175 tion can take at least one of the following roles:
 176 **provisioning** if it delivers local event logs to
 177 be collaboratively mined; **mining** if it applies

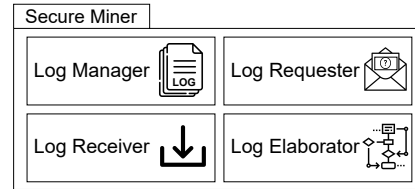


Figure 3: Sub-components of the Secure Miner

178 process mining algorithms using event logs retrieved from provisioners. In Fig. 2,
 179 we propose the high-level schematization of the CONFINE framework. In our
 180 solution, every organization hosts one or more nodes hosting components (the
 181 names of which will henceforth be formatted with a teletype font). Depending
 182 on the played role, nodes come endowed with a Provisioner or a Secure Miner
 183 component, or both. The Provisioner component consists of the following two
 184 main sub-components. The Log Recorder registers the events taking place in the
 185 organizations' systems. The Log Provider delivers on-demand data to mining
 186 players. The hospital and all other parties in our example record Alice and Bob's
 187 cases using the Log Recorder. The Log Recorder is queried by the Log Provider
 188 for event logs to be made available for mining. The latter controls access to local

189 event logs by authenticating data requests by miners and rejecting those that come
190 from unauthorized parties. In our motivating scenario, the specialized clinic, the
191 pharmaceutical company, and the hospital leverage Log Providers to authenticate
192 the miner party before sending their logs. The Secure Miner component shelters
193 external event logs inside a protected environment to preserve data confidentiality
194 and integrity. Notice that Log Providers accept requests issued solely by Secure
195 Miners. Next, we provide an in-depth focus on the latter.

196 **The Secure Miner.** The primary objective of the Secure Miner is to allow
197 miners to securely execute process mining algorithms using event logs retrieved
198 from provisioners such as the specialized clinic, pharmaceutical company, and the
199 hospital in our example. Secure Miners are isolated components that guarantee
200 data inalterability and confidentiality. In [Fig. 3](#), we show a schematization of
201 the Secure Miner, which consists of four sub-components: (i) Log Requester;
202 (ii) Log Receiver; (iii) Log Manager; (iv) Log Elaborator. The Log Requester
203 and the Log Receiver are the sub-components that we employ during the event log
204 retrieval. Log Requesters send authenticable data requests to the Log Providers.
205 The Log Receiver collects event logs sent by Log Providers and entrusts them
206 to the Log Manager, securing them from accesses that are external to the Secure
207 Miner. Miners of our motivating scenario, such as the university and the national
208 institute of statistics, employ these three components to retrieve and store Alice
209 and Bob’s data. The Log Elaborator merges the event data locked in the Secure
210 Miner to have a global view of the inter-organizational process comprehensive of
211 activities executed by each involved party. Thereupon, it executes process mining
212 algorithms in a protected environment, inaccessible from the outside computation
213 environment. In our motivating scenario, the Log Elaborator combines the traces

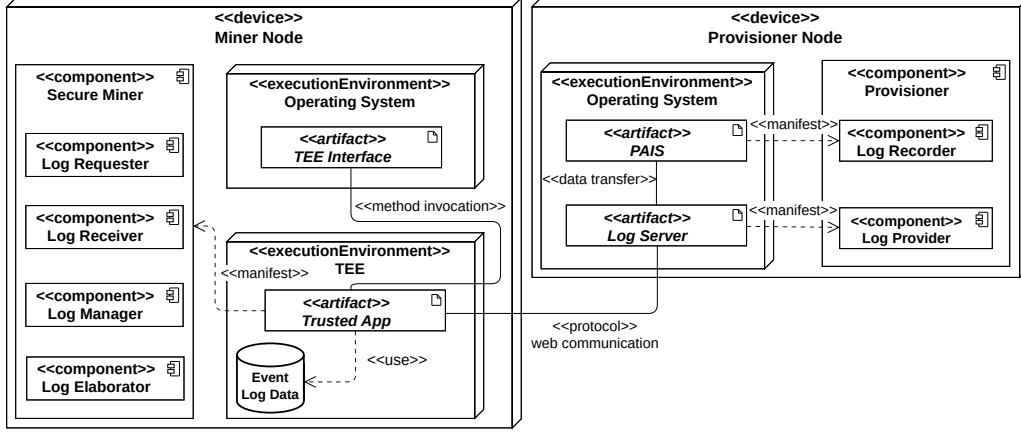


Figure 4: UML deployment diagram of the CONFINE architecture

214 associated with the cases of Alice (i.e., T_{312}^H , T_{312}^S , and T_{312}^C) and Bob (i.e., T_{711}^H , T_{711}^S ,
 215 and T_{711}^C), generates the chronologically sorted traces T_{312} and T_{711} , and feeds them
 216 into the mining algorithms (see the bottom-right quadrant of Sect. 2.2).

217 6. Realization

218 In this section, we outline the technical aspects concerning the realization of
 219 our solution. Therefore, we first present the enabler technologies through which
 220 we instantiate the design principles presented in Sect. 5. After that, we discuss the
 221 CONFINE interaction protocol. Finally, we show the implementation details.

222 6.1. Deployment

223 Figure 4 depicts a UML deployment diagram [15] to illustrate the employed tech-
 224 nologies and computation environments. We recall that the Miner and Provisioner
 225 nodes are drawn as separated, although organizations can host both. In our motivating
 226 scenario, e.g., the hospital can be equipped with machines aimed for both mining
 227 and provisioning.

228 Provisioner Nodes host the Provisioner’s components, encompassing the
229 Log Recorder and the Log Provider. The Process-Aware Information System
230 (PAIS) manifests the Log Recorder [16]. The PAIS grants access to the Log
231 Server, enabling it to retrieve event log data. The Log Server, on the other hand,
232 embodies the functionalities of the Log Provider, implementing web services
233 aimed at handling remote data requests and providing event log data to miners.

234 The Miner Node is characterized by two distinct *execution environments*: the
235 Operating System (OS) and the Trusted Execution Environment (TEE) [5]. TEEs
236 establish isolated contexts separate from the OS, safeguarding code and data through
237 hardware-based encryption mechanisms. This technology relies on dedicated
238 sections of a CPU capable of handling encrypted data within a reserved section
239 of the main memory [17]. By enforcing memory access restrictions, TEEs aim to
240 prevent one application from reading or altering the memory space of another, thus
241 enhancing the overall security of the system. This dedicated areas in memory are,
242 however, limited. Once the limits are exceeded, TEEs have to scout around in outer
243 memory areas, thus conceding the opportunity to malicious reader to understand
244 the saved data based on the memory reads and writes. To avoid this risk, TEE
245 implementations often raise errors that halt the program execution when the memory
246 demand goes beyond the available space. Therefore, the design of secure systems
247 that resort to TEEs must take into account that memory consumption must be kept
248 under control. We leverage the security guarantees provided by TEEs [18] to protect
249 a Trusted App responsible for fulfilling the functions of the Secure Miner and
250 its associated sub-components. The TEE ensures the integrity of the Trusted App
251 code, protecting it against potential malicious manipulations and unauthorized
252 access by programs running within the Operating System. Additionally, we utilize

the isolated environment of TEEs to securely store event log data (e.g., Alice and Bob’s cases). The TEE retains a private key in the externally inaccessible section of memory, paired with a public key in a Rivest-Shamir-Adleman (RSA) [19] scheme for attestation (only the owner of the private key can sign messages that are verifiable via the public key) and secure message encryption (only the owner of the private key can decode messages that are encrypted with the corresponding public key). In our solution, access to data located in the TEE is restricted solely to the Trusted App. Users interact with the Trusted App through the Trusted App Interface, which serves as the exclusive communication channel. The Trusted App offers secure methods, invoked by the Trusted App Interface, for safely receiving information from the Operating System and outsourcing the results of computations.

6.2. The CONFINE protocol

We orchestrate the interaction of the components in CONFINE via a protocol. We separate it in four subsequent stages, namely (i) *initialization*, (ii) *remote attestation*, (iii) *data transmission*, and (iv) *computation*. These stages are depicted in Figs. 5(a), 5(b), 6(a) and 6(b), respectively. Our protocol involves two primary entities: a Secure Miner (hereafter referred to as \mathcal{M}) and one or more Provisioners ($\mathcal{P}_1, \dots, \mathcal{P}_n \in \hat{\mathcal{P}}$). The behavioral descriptions for \mathcal{M} and \mathcal{P}_i are outlined in Alg. 1 and Alg. 2, respectively. These specifications adhere to the distributed programming syntax detailed in [20]. We assume that communication between Secure Miners and Log Provisioners occurs through an *Authenticated Point-to-Point Perfect Link* [20]. This communication abstraction guarantees: (i) *reliable delivery*, (ii) *no duplication*, and (iii) *authenticity*. In order to enhance clarity, we have adapted the foundational notations of the SEND and DELIVER events of this primitive to emphasize message senders (indicated by the symbol ‘ \ll ’ when DELIVER occurs) and receivers

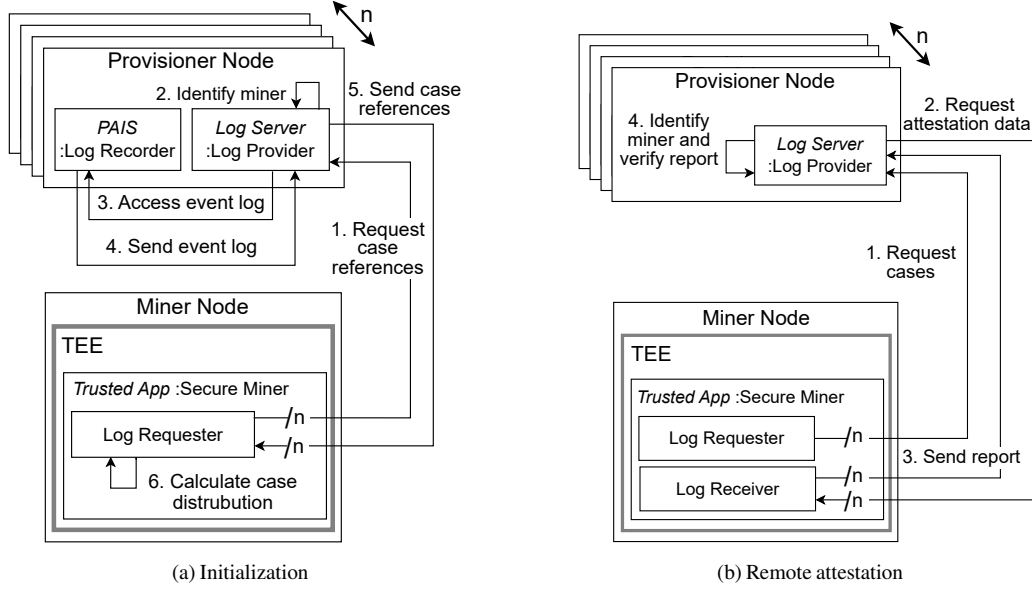


Figure 5: Unfolding example for the initialization, remote attestation phases of the CONFINE protocol

278 (indicated by ‘ \gg ’ when SEND is triggered). In Alg. 1, the Secure Miner is provided
 279 with the following inputs: the list of Provisioners’ references ($\mathcal{P}_1, \dots, \mathcal{P}_n$) and a
 280 segment size *seg_size* employed for the log segmentation during the *data trasmission*
 281 phase. Similarly, the Provisioner’s specification in Alg. 2 considers as input the
 282 list of references to miners ($\mathcal{M}_1, \dots, \mathcal{M}_s$) for which event log access is enabled.

283 In the following, we describe each protocol phase in detail.

284 **Initialization.** The objective of the initialization stage is to inform the miner about the
 285 distribution of cases related to a business process among the Provisioner Nodes.
 286 At the onset of this stage, the Log Requester within the Trusted App issues n
 287 requests, one per Log Server component, to retrieve the list of case references they
 288 record (step 1 in Fig. 5(a) and Alg. 1, line 3). Following sender authentication (2),
 289 each Log Server retrieves the local event log from the PAIS (3, 4) and subsequently
 290 responds to the Log Requester by providing a list of its associated case references

291 (5 and Alg. 2, line 3). After collecting these n responses (Alg. 1, line 4), the Log
 292 Requester delineates the distribution of cases. In the context of our motivating
 293 scenario, by the conclusion of the initialization, the miner gains knowledge that
 294 the case associated with Bob, synthesized in the traces T_{711}^H and T_{711}^C , is exclusively
 295 retained by the hospital and the specialized clinic. In contrast, the traces of Alice's
 296 case, denoted as T_{312}^H , T_{312}^C , and T_{312}^S , are scattered across all three organizations.
 297 **Remote attestation.** The remote attestation serves the purpose of establishing trust
 298 between miners and provisioners in the context of fulfilling data requests. This phase
 299 adheres to the overarching principles outlined in the RATS RFC standard [21] serving
 300 as the foundation for several TEE attestation schemes (e.g., Intel EPID,² and AMD
 301 SEV-SNP³). Remote attestation has a dual objective: (i) to furnish provisioners with
 302 compelling evidence that the data request for an event log originates from a Trusted
 303 App running within a TEE; (ii) to confirm the specific nature of the Trusted App
 304 as an authentic Secure Miner software entity. This phase is triggered when the
 305 Log Requester sends a new case request to the Log Server (step 1 in Fig. 5(b) and
 306 Alg. 2, line 5), specifying: (i) the segment size (henceforth, seg_size), and (ii) the
 307 set of the requested case *IIDs*. Both parameters will be used in the subsequent *data*
 308 *transmission* phase. Each of the n Log Servers commences the verification process
 309 by requesting the necessary information from the Log Receiver to conduct the
 310 attestation (2). Subsequently, the Log Receiver generates the attestation report
 311 containing the so-called *measurement* of the Trusted App, which is defined as
 312 the hash value of the combination of its source code and data. Once this report is

²sgx101.gitbook.io/sgx101/sgx-bootstrap/attestation. Accessed: 17/01/2024.

³amd.com/en/processors/amd-secure-encrypted-virtualization. Accessed: 17/01/2024.

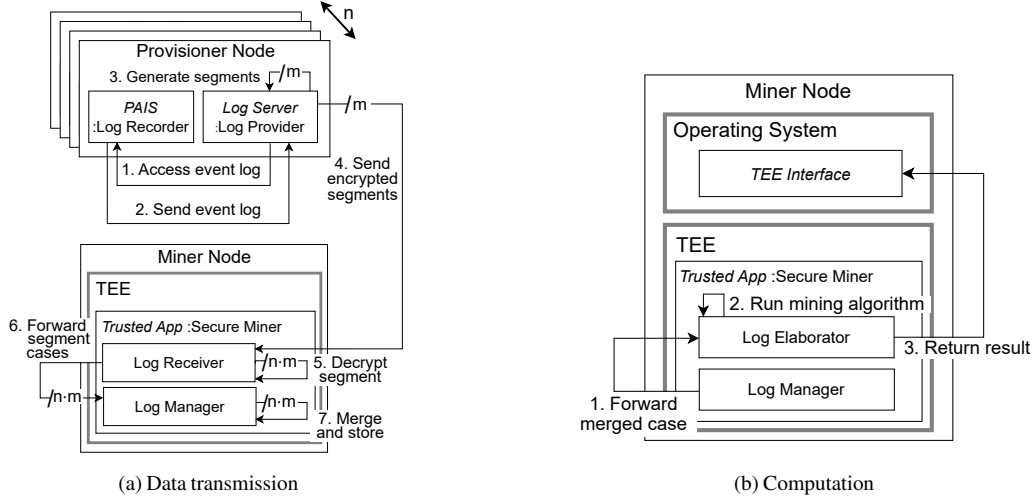


Figure 6: Unfolding example for the data transmission and computation phases of the CONFINE protocol

313 signed using the attestation private key associated with the TEE’s hardware of the
 314 Miner Node, it is transmitted by the Log Receiver to the Log Servers alongside
 315 the attestation public key of the Miner Node (3). The Log Servers authenticate
 316 the miner using the public key and decrypt the report (4). In this last step, the Log
 317 Servers undertake a comparison procedure in which they juxtapose the measurement
 318 found within the decrypted report against a predefined reference value associated
 319 with the source code of the Secure Miner. If the decrypted measurement matches
 320 the predefined value, the Miner Node gains trust from the provisioner.

321 **Data transmission.** Once the trusted nature of the Trusted App is verified, the Log
 322 Servers proceed with the transmission of their cases. To accomplish this, each Log
 323 Server retrieves the event log from the PAIS (steps 1 and 2 in Fig. 6(a)), and filters
 324 it according to the case reference set specified by the miner. Given the constrained
 325 workload capacity of the TEE, it is imperative for Log Servers to partition the filtered
 326 event log into distinct segments. Consequently, each Log Server generates m log

327 segments comprising a variable count of entire cases (3 and Alg. 2, line 6). The
 328 cumulative size of these segments is governed by the threshold parameter specified
 329 by the miner in the initial request (step 1 of the remote attestation phase, Fig. 5(b)).
 330 As an illustrative example from our motivating scenario, the Log Server of the
 331 hospital may structure the segmentation such that T_{312}^H and T_{711}^H reside within the
 332 same segment, whereas the specialized clinic might have T_{312}^S and T_{711}^S in separate
 333 segments. Subsequently, the n Log Servers transmit their m encrypted segments to
 334 the Log Receiver of the Trusted App (4 and Alg. 2, line 8). The Log Receiver,
 335 in turn, collects the $n \times m$ responses in a queue, processing them one at a time.
 336 After decrypting the processed segment (5), the Log Receiver forwards the cases
 337 contained in the segment to the Log Manager (6 and Alg. 1, line 15). To reconstruct
 338 the process instance, cases belonging to the same process instance must be merged by
 339 the Log Manager resulting in a single trace (e.g., T_{312} for Alice case) comprehensive
 340 of all the events in the partial traces (e.g., T_{312}^H , T_{312}^S and T_{312}^C for Alice case). During
 341 this operation, the Log Manager applies a specific *merging schema* (i.e., a rule
 342 specifying the attributes that link two cases during the merge) as stated in [12]. In our
 343 illustrative scenario, the merging schema to combine the cases of Alice is contingent
 344 upon the linkage established through their case identifier (i.e., 312). We underline
 345 that our proposed solution facilitates the incorporation of diverse merging schemas
 346 encompassing distinct trace attributes. The outcomes arising from merging the cases
 347 within the processed segment are securely stored by the Log Manager in the TEE.
 348 **Computation.** The Trusted App requires all the provisioners to have delivered
 349 cases referring to the same process instances. For example, when the hospital and
 350 the other organizations have all delivered their information concerning case 312
 351 to the Trusted App, the process instance associated with Alice becomes eligible

352 for computation. Upon meeting this condition (Alg. 1, line 16), the Log Manager
353 forwards the case earmarked for computation to the Log Elaborator (step 1 in
354 Fig. 6(b) and Alg. 1, line 18). Subsequently, the Log Elaborator proceeds to input
355 the merged case into the process mining algorithm (2). Ultimately, the outcome
356 of the computation is relayed by the Log Elaborator from the TEE to the TEE
357 Interface running atop the Operating System of the Miner Node (3). The
358 CONFINE protocol does not impose restrictions on the post-computational handling
359 of results. In our motivating scenario, the University and the National Institute of
360 Statistics, serving as miners, disseminate the outcomes of computations, generating
361 analyses that benefit the provisioners (though the original data are never revealed in
362 clear). Furthermore, our protocol enables the potential for provisioners to have their
363 proprietary Secure Miner, allowing them autonomous control over the computed
364 results.

365 6.3. Implementation

366 We implemented the Secure Miner component as an Intel SGX⁴ trusted ap-
367 plication, encoded in Go through the EGo framework.⁵ We resort to a TLS [22]
368 communication channel between miners and provisioners over the HTTP web pro-
369 tocol to secure the information exchange. To demonstrate the effectiveness of our
370 framework, we re-implemented and integrated the *HeuristicsMiner* discovery algo-
371 rithm [23] within the Trusted Application. Our implementation of CONFINE,
372 including the *HeuristicsMiner* in Go, is openly accessible at the following URL:
373 github.com/Process-in-Chains/CONFINE/.

⁴sgx101.gitbook.io/sgx101/. Accessed: 17/01/2024.

⁵docs.edgeless.systems/ego. Accessed: 17/01/2024.

Algorithm 1: Secure Miner's behavior in CONFINE.

Input: $\hat{\mathcal{P}} = \{\mathcal{P}_1, \dots, \mathcal{P}_n\}$, the (references to) n log provisioners;
 seg_size , the maximum size of the log segment to be transmitted by the log provisioners.

Data: $IIDMap: \widehat{IID} \rightarrow 2^{\hat{\mathcal{P}}}$, a map from case references $iid \in \widehat{IID}$ to the set of log provisioners in $\hat{\mathcal{P}}$;
 $PMap: \hat{\mathcal{P}} \rightarrow 2^{\widehat{IID}}$, a map from log provisioners $\mathcal{P} \in \hat{\mathcal{P}}$ to the set of references to their cases in \widehat{IID} ;
 $Cases: \widehat{IID} \rightarrow \hat{\mathcal{C}}$, a map from case references $iid \in \widehat{IID}$ to a set of cases in $\hat{\mathcal{C}}$.

Implements: SecureMiner, instance \mathcal{M} .

Uses: AuthenticatedPerfectPointToPointLink, instance aL .

```

1  upon event  $\langle \mathcal{M}, \text{INIT} | \hat{\mathcal{P}}, seg\_size \rangle$  do           // The Log Requester of  $\mathcal{M}$  starts the CONFINE protocol – initialization phase in Fig. 5(a)
2    foreach  $\mathcal{P} \in \hat{\mathcal{P}}$  do                               // For every Provisioner  $\mathcal{P}$ 
3      trigger  $\langle aL, \text{SEND} \gg \mathcal{P} | \text{CASESREFREQ} \rangle$       // Request  $\mathcal{P}$ 's case references (see Alg. 2, line 1)
4    upon  $|\hat{\mathcal{P}}| = |\text{dom}(PMap)|$  do                       // Once all Provisioners have answered with their case references
5      foreach  $\mathcal{P} \in \hat{\mathcal{P}}$  do trigger  $\langle aL, \text{SEND} \gg \mathcal{P} | \text{CASESREQ}, seg\_size, PMap[\mathcal{P}] \rangle$  // Request their cases via  $aL$  (see Alg. 2, line 4)

6  upon event  $\langle aL, \text{DELIVER} \ll \mathcal{P} | \text{CASESREFRES}, IIDs \rangle$  such that  $\mathcal{P} \in \hat{\mathcal{P}}$  do //  $\mathcal{M}$ 's Log Requester gets  $\mathcal{P}$ 's case references via  $aL$  (Alg. 2, line 3)
7    foreach  $iid \in IIDs$  do                             // For every case reference  $iid$  received in  $IIDs$ 
8       $IIDMap[iid] \leftarrow IIDMap[iid] \cup \{\mathcal{P}\}$       // Add  $\mathcal{P}$  to the set of provisioners for case  $iid$  in  $IIDMap$ 
9       $PMap[\mathcal{P}] \leftarrow PMap[\mathcal{P}] \cup IIDs$            // Register the references of the cases provided by  $\mathcal{P}$  in  $PMap$ 

10 upon event  $\langle aL, \text{DELIVER} \ll \mathcal{P}, [\text{CASESRES}, S] \rangle$  such that  $\mathcal{P} \in \hat{\mathcal{P}}$  do //  $\mathcal{M}$ 's Log Receiver gets a segment from  $\mathcal{P}$  via  $aL$  (Alg. 2, line 8)
11   foreach  $C_{iid}^{\mathcal{P}} \in S$  do                             // For every  $C_{iid}^{\mathcal{P}}$  in the delivered segment  $S$ , each associated with a  $iid$  – data transmission phase in Fig. 6(a)
12     if  $iid \in PMap[\mathcal{P}]$  then                             // If  $\mathcal{P}$  has declared the ownership of  $iid$  (see line 6)
13        $PMap[\mathcal{P}] \leftarrow PMap[\mathcal{P}] \setminus \{iid\}$       // Remove  $iid$  from the set of case references to be provided by  $\mathcal{P}$ 
14        $IIDMap[iid] \leftarrow IIDMap[iid] \setminus \{\mathcal{P}\}$  // Remove  $\mathcal{P}$  from the set of  $iid$  provisioners
15        $\mathcal{M}.\text{LogManager.mergeAndStore}(Cases, C_{iid}^{\mathcal{P}})$  // Update the case via  $\oplus$  and store the result in  $Cases$ 

16 upon  $IIDMap[iid] = \emptyset$  for some  $iid \in \text{dom}(IIDMap)$  do // When all the pieces of some  $iid$  have arrived to  $\mathcal{M}$ 's Log Manager
17    $\text{dom}(IIDMap) \leftarrow \text{dom}(IIDMap) \setminus \{iid\}$  // Remove  $iid$  from the domain of cases which still needs to be processed
18   yield  $Cases[iid]$  to  $\mathcal{M}.\text{LogElaborator}$  // Forward the case  $iid$  to the Log Elaborator of  $\mathcal{M}$  for mining – computation phase in Fig. 6(b)

```

Algorithm 2: Provisioner's behavior in CONFINE.

Input: $\hat{\mathcal{M}} = \{\mathcal{M}_1, \dots, \mathcal{M}_s\}$, the (references to) s miners.

Implements: Provisioner, instance \mathcal{P} .

Uses: AuthenticatedPerfectPointToPointLink, instance aL .

```

1  upon event  $\langle aL, \text{DELIVER} \ll \mathcal{M} | \text{CASESREFSREQ} \rangle$  such that  $\mathcal{M} \in \hat{\mathcal{M}}$  do //  $\mathcal{P}$  receives the request for case references from  $\mathcal{M}$  (see Alg. 1, line 3)
2     $IIDs \leftarrow \mathcal{P}.\text{LogRecorder.accessCaseReferences}()$  // Access the case references via Log Recorder
3    trigger  $\langle aL, \text{SEND} \gg \mathcal{M} | \text{CASESREFRES}, IIDs \rangle$  // send the case references to  $\mathcal{M}$  (see Alg. 1, line 6)

4  upon event  $\langle aL, \text{DELIVER} \ll \mathcal{M} | \text{CASESREQ}, seg\_size, IIDs \rangle$  such that  $\mathcal{M} \in \hat{\mathcal{M}}$  do //  $\mathcal{P}$  gets the case request from  $\mathcal{M}$  (see Alg. 1, line 5)
5    if  $\mathcal{M}.\text{LogReceiver.getAttestationReport}(\mathcal{P})$  is valid then // Get and verify the attestation report of  $\mathcal{M}$  – remote attestation in Fig. 5(b)
6       $\{S_1, \dots, S_m\} \leftarrow \mathcal{P}.\text{LogProvider.segmentEventLog}(\mathcal{P}.\text{LogRecorder.accessEventLog}(IIDs), seg\_size)$  // Segment the event log
7      foreach  $i \in \{1, \dots, m\}$  do                         // For every split segment  $S_i$ 
8        trigger  $\langle aL, \text{SEND} \gg \mathcal{M} | \text{CASESRES}, S_i \rangle$  // send the segment  $S_i$  to  $\mathcal{M}$  (see Alg. 1, line 10) – data transmission phase in Fig. 6(a)

```

7. Evaluation

In this section, we evaluate our approach through the testing of our tool implementation. We begin with a convergence analysis to demonstrate the correctness of the collaborative data exchange process. Subsequently, we gauge the memory usage with synthetic and real-life event logs, to observe the trend during the enactment of our

Table 2: Event logs used for our experiments

| Name | Type | Activities | Cases | Max events | Min events | Avg. events | Organization \mapsto Activities |
|---------------------|-----------|------------|-------|------------|------------|-------------|--|
| Motivating scenario | Synthetic | 19 | 1000 | 18 | 9 | 14 | $\mathcal{O}^P \mapsto 3, \mathcal{O}^C \mapsto 5, \mathcal{O}^H \mapsto 14$ |
| Sepsis [24] | Real | 16 | 1050 | 185 | 3 | 15 | $\mathcal{O}^1 \mapsto 1, \mathcal{O}^2 \mapsto 1, \mathcal{O}^3 \mapsto 14$ |
| BPIC2013 [25] | Real | 7 | 1487 | 123 | 1 | 9 | $\mathcal{O}^1 \mapsto 6, \mathcal{O}^2 \mapsto 7, \mathcal{O}^3 \mapsto 6$ |

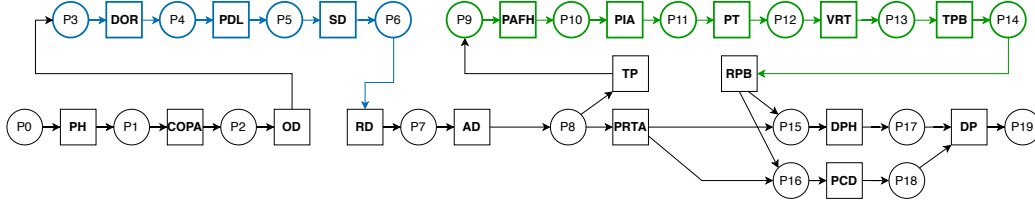


Figure 7: *HeuristicsMiner* output in CONFINE

379 protocol and assess scalability. We recall that we focus on memory utilization since
 380 the availability of space in the dedicated areas is limited as we discussed in Sect. 6.1.
 381 We discuss our experimental results in the following. For the sake of reproducibility,
 382 we make available all the testbeds and results in our public code repository (linked
 383 above).

384 **Output convergence.** To experimentally validate the correctness of our approach
 385 in the transmission and computation phases (see Sect. 6), we run a *convergence*
 386 test. To this end, we created a synthetic event log consisting of 1000 cases of 14
 387 events on average (see Table 2) by simulating the inter-organizational process of
 388 our motivating scenario (see Fig. 1)⁶ and we partitioned it in three sub-logs (one
 389 per involved organization), an excerpt of which is listed in Sect. 2.2. We run the
 390 stand-alone *HeuristicsMiner* on the former, and processed the latter through our
 391 CONFINE toolchain. As expected, the results converge and are depicted in Fig. 7 in

⁶We generated the event log through BIMP (<https://bimp.cs.ut.ee/>). We filtered the generated log by keeping the sole events that report on the completion of activities, and removing the start and end events of the pharmaceutical company and specialized clinic’s sub-processes.

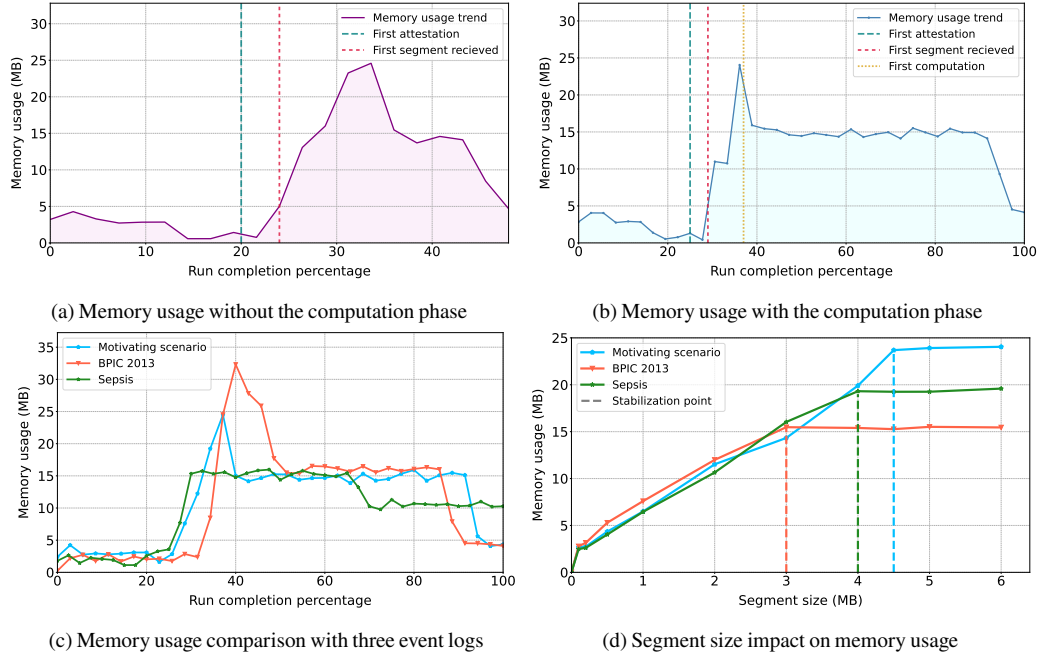


Figure 8: Memory usage test results

the form of a workflow net [26]. For clarity, we have colored activities recorded by the organizations following the scheme of Table 2 (black for the hospital, blue for the pharmaceutical company, and green for the specialized clinic).

Memory usage. Figures 8(a) and 8(b) display plots corresponding to the runtime memory utilization of our CONFINE implementation (in MegaBytes). Differently from Fig. 8(b), Fig. 8(a) excludes the computation stage by leaving the *HeuristicsMiner* inactive so as to isolate the execution from the mining-specific operations. The dashed lines mark the starting points for the remote attestation, the data transmission and the computation stages. We held the *seg_size* constant at 2000 KiloBytes. We observe that the data transmission stage reaches the highest peak of memory utilization, which is then partially freed by the subsequent computation stage, steadily occupying memory space at a lower level. To verify whether this

phenomenon is due to the synthetic nature of our simulation-based event log, we also gauge the runtime memory usage of two public real-world event logs too (Sepsis [24] and BPIC 2013 [25]). The characteristics of the event logs are summarized in Table 2. Since those are *intra-organizational* event logs, we split the contents to mimic an *inter-organizational* context. In particular, we separated the Sepsis event log based on the distinction between normal-care and intensive-care paths, as if they were conducted by two distinct organizations. Similarly, we processed the BPIC 2013 event log to sort it out into the three departments of the Volvo IT incident management system. Figure 8(c) depicts the results. We observe that the processing of the BPIC 2013 event log demands more memory, particularly during the initial stages, probably owing to its larger size. Conversely, the Sepsis event log turns out to entail the least expensive run. To verify whether these trends are affected by the dimension of the exchanged data segments, we conducted an additional test to examine the trend of memory usage as the *seg_size* varies with all the aforementioned event logs. Notably, the polylines displayed in Fig. 8(d) indicate a linear increment of memory occupation until a breakpoint is reached. After that, the memory in use is steady. These points, marked by vertical dashed lines, correspond to the *seg_size* value that allows the provider's segments to be contained in a single data segment.

Scalability. In this subsection, we examine the scalability of the Secure Miner, focusing on its capacity to efficiently manage an increasing workload in the presence of limited memory resources. We implemented three distinct test configurations gauging runtime memory usage as variations of our motivating scenario log. In particular, we considered (I) the maximum number of events per case, (II) the number of cases $|\widehat{IID}|$, and (III) the number of provisioning organizations $|\widehat{O}|$ as independent integer variables. To conduct the test on the maximum number of

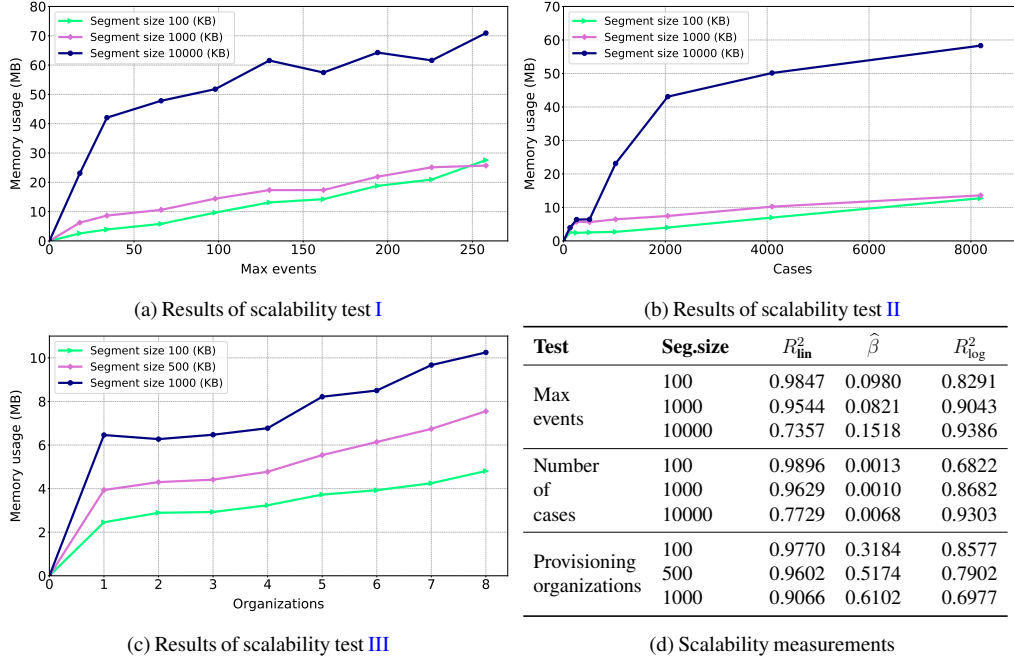


Figure 9: Scalability test results

429 events, we added a loop back from the final to the initial activity of the process model,
 430 progressively increasing the number of iterations $2 \leq x_{\mathcal{G}} \leq 16$ at a step of 2, resulting
 431 in $18 + 16 \cdot (x_{\mathcal{G}} - 1)$ events. Concerning the test on the number of cases, we simulated
 432 additional process instances so that $|\widehat{IID}| = 2^{x_{iid}}$ having $x_{iid} \in \{7, 8, \dots, 13\}$. Finally, for
 433 the assessment of the number of organizations, the test necessitated the distribution
 434 of the process model activities' into a variable number of pools, each representing a
 435 different organization ($|\hat{\mathcal{O}}| \in \{1, 2, \dots, 8\}$). We parameterized the above configurations
 436 with three segment sizes (in KiloBytes): $seg_size \in \{100, 1000, 10000\}$ for tests
 437 I and II, and $seg_size \in \{100, 500, 1000\}$ for test III (the range is reduced without loss
 438 of generality to compensate the partitioning of activities into multiple organizations).
 439 To facilitate a more rigorous interpretation of the output trends across varying
 440 seg_size configurations, we employ two well-known statistical measures. As a

primary measure of goodness-of-fit, we employ the coefficient of determination R^2 [27], which assesses the degree to which the observed data adheres to the linear (R_{lin}^2) and logarithmic (R_{log}^2) regressions derived from curve fitting approximations. To further delve into the analysis of trends with a high R_{lin}^2 , we consider the slope $\hat{\beta}$ of the approximated linear regression [28].

Table 9(d) lists the measurements we obtained. We describe them to elucidate the observed patterns. Figure 9(a) depicts the results of test I, focusing on the increase of memory utilization when the number of events in the event logs grows. We observe that the memory usage for *seg_size* 100 and 1000 (depicted by green and lilac lines, respectively) are quite similar, whereas the setting with *seg_size* 10,000 (blue line) exhibits significantly higher memory usage. For the settings with *seg_size* 100 and 1000, R_{lin}^2 approaches 1, signifying an almost perfect approximation of the linear relation, against lower R_{log}^2 values. In these test settings, $\hat{\beta}$ is very low yet higher than 0, thus indicating that memory usage is likely to continue increasing as the number of max events grows. The configuration with *seg_size* 10,000 yields a higher R_{log}^2 value, thus suggesting a logarithmic trend, hence a greater likelihood of stabilizing memory usage growth rate as the number of maximum events increases. In Fig. 9(b), we present the results of test II, assessing the impact of the number of cases on the memory consumption. As expected, the configurations with *seg_size* set to 100 and 1000 exhibit a trend of lower memory usage than settings with *seg_size* 10,000. The R_{lin}^2 score of the trends with *seg_size* 100 and 1000 indicate a strong linear relationship between the dependent and independent variables compared to the trend with *seg_size* 10,000, which is better described by a logarithmic regression ($R_{\text{log}}^2 = 0.9303$). For the latter, the R_{log}^2 value is higher than the corresponding R_{lin}^2 thus suggesting that the logarithmic approximation is better suited to describe the

466 trend. Differently from test I, the $\hat{\beta}$ score associated with the linear approximations
 467 of the trends with *seg_size* 100 and 1000 approaches 0, indicating that the growth
 468 rate of memory usage as the number of cases increases is negligible. In Fig. 9(c), we
 469 present the results of test III, on the relation between the number of organizations
 470 and the memory usage. The chart shows that memory usage trends increase as
 471 provisioning organizations increase for all three segment sizes. The R_{lin}^2 values
 472 for the three *seg_sizes* are very high, indicating a strong positive linear correlation.
 473 The test with *seg_size* 100 exhibits the slowest growth rate, as corroborated by the
 474 lowest $\hat{\beta}$ result (0.3184). For the configuration with *seg_size* 500, the memory usage
 475 increases slightly faster ($\hat{\beta} = 0.5174$). With *seg_size* 1000, the overall memory usage
 476 increases significantly faster than the previous configurations ($\hat{\beta} = 0.6102$). We
 477 derive from these findings that the Secure Miner may encounter scalability issues
 478 when handling settings with a large number of provisioning organizations. Further
 479 investigation is warranted to determine the precise cause of this behavior and identify
 480 potential mitigation strategies.

481 In the next section, we conclude our work and outline future research directions
 482 based upon our current findings and the limitations of our approach.

483 **8. Conclusion and Future Work**

484 Confidentiality is paramount in inter-organizational process mining due to the
 485 transmission of sensitive data across organizational boundaries. Our research
 486 investigates a decentralized secrecy-preserving approach that enables organizations
 487 to employ process mining techniques with event logs from multiple organizations
 488 while ensuring the protection of privacy and confidentiality. Our solution offers
 489 a number of directions to walk along for improvement. We operate under the

490 assumption of fair conduct by data provisioners and do not account for the presence
491 of injected or maliciously manipulated event logs. In addition, we assume that miners
492 and provisioners exchange messages in reliable communication channels where no
493 loss or bit corruption occurs. Our approach relies on certain assumptions about event
494 log data, including the existence of a universal clock for event timestamps, which may
495 not be realistic in situations where organizations are not perfectly synchronized. We
496 aim at enhancing our approach to make it robust to the relaxation of these constraints.
497 Our future work encompasses the integration of usage control policies that specify
498 rules on event logs' utilization. We plan to design policy enforcement and monitoring
499 mechanisms to achieve this goal following the principles already addressed in [29, 30].
500 Our solution embraces process mining techniques in a general way. However, we
501 believe the presented approach is compatible with declarative model representations
502 [31]. Therefore, trusted applications could compute and store the entire set of
503 rules representing a business process, and users may interact with them via trusted
504 queries. Finally, in our implementation, we have focused on process discovery tasks.
505 Nevertheless, our approach has the potential to seamlessly cover a wider array of
506 process mining functionalities such as *conformance checking*, and *performance*
507 *analysis* techniques. Implementing them and showing their integrability with our
508 approach paves the path for future research endeavors.

509 **Acknowledgments.** The authors thank Giuseppe Ateniese for the fruitful discussion
510 and insights. This research work was partly funded by MUR under PRIN grant
511 B87G22000450001 (PINPOINT), by the Latium Region under PO FSE+ grant
512 B83C22004050009 (PPMPP), and by the EU-NGEU under the NRRP MUR grant
513 PE00000014 (SERICS).

EXTENSION PLAN:

Extended introduction

Add Background Section

Add more related work

Modify the motivating scenario (and the design alongside the implementation) with more attributes involved in the event log (for example, the famous ID that might not be shared among providers, many of those having differing attributes and codes to refer to an instance).

Add Notation and formalization of event logs, merging, partitioning and segmentation

Add Soundness and completeness theorems

Add threat model

Full pseudocode of the protocol ✓

More real-world event logs (plus two, at least) with associated tests

Add communication overhead v. segment size charts: elab time vs segment size; total time (incl. network) vs segment size.

Integrate declarative conformance checking (let it be with Janus or MINERful).

514

References

515

516 [1] W. M. P. van der Aalst, et al., Process mining manifesto, in: BPM Workshops,
517 2012, pp. 169–194.

518 [2] W. M. P. van der Aalst, Intra-and inter-organizational process mining: Dis-
519 covering processes within and between organizations, in: PoEM, 2011, pp.
520 1–11.

521 [3] C. Liu, Q. Li, X. Zhao, Challenges and opportunities in collaborative business
522 process management: Overview of recent advances and introduction to the
523 special issue, Inf. Syst. Front. 11 (2009) 201–209.

524 [4] M. Müller, N. Ostern, Koljada, et al., Trust mining: analyzing trust in
525 collaborative business processes, IEEE Access (2021) 65044–65065.

526 [5] M. Sabt, M. Achemlal, A. Bouabdallah, Trusted execution environment: What

- 527 it is, and what it is not, in: 2015 IEEE TrustCom/BigDataSE/ISPA, 2015, pp.
528 57–64.
- 529 [6] M. Müller, A. Simonet-Boulogne, S. Sengupta, O. Beige, Process mining in
530 trusted execution environments: Towards hardware guarantees for trust-aware
531 inter-organizational process analysis, in: ICPM, 2021, pp. 369–381.
- 532 [7] G. Elkoumy, S. A. Fahrenkrog-Petersen, et al., Shareprom: A tool for privacy-
533 preserving inter-organizational process mining, in: BPM (PhD/Demos), 2020,
534 pp. 72–76.
- 535 [8] G. Elkoumy, S. A. Fahrenkrog-Petersen, et al., Secure multi-party computation
536 for inter-organizational process mining, in: BPMDS/EMMSAD, 2020, pp.
537 166–181.
- 538 [9] W. M. P. van der Aalst, Federated process mining: Exploiting event data across
539 organizational boundaries, in: SMDS 2021, 2021, pp. 1–7.
- 540 [10] R. Cramer, I. Damgård, J. B. Nielsen, Secure Multiparty Computation and
541 Secret Sharing, Cambridge University Press, 2015.
- 542 [11] C. Zhao, S. Zhao, Zhao, et al., Secure multi-party computation: Theory,
543 practice and applications, Inf. Sci. 476 (2019) 357–372.
- 544 [12] J. Claes, G. Poels, Merging event logs for process mining: A rule based
545 merging method and rule suggestion algorithm, Expert Syst. Appl. 41 (2014)
546 7291–7306.
- 547 [13] J. D. Hernandez-Resendiz, E. Tello-Leal, H. M. Marin-Castro, et al., Merging
548 event logs for inter-organizational process mining, in: New Perspectives

- 549 on Enterprise Decision-Making Applying Artificial Intelligence Techniques,
550 Springer, 2021, pp. 3–26.
- 551 [14] M. Jans, M. Hosseinpour, How active learning and process mining can act as
552 continuous auditing catalyst, *Int. J. Accounting Inf. Systems* 32 (2019) 44–58.
- 553 [15] N. Koch, A. Kraus, The expressive power of UML-based web engineering, in:
554 IWWOST02, volume 16, 2002, pp. 40–41.
- 555 [16] M. Dumas, M. La Rosa, J. Mendling, H. A. Reijers, *Fundamentals of Business
556 Process Management, Second Edition*, Springer, 2018.
- 557 [17] V. Costan, S. Devadas, Intel SGX explained, *Cryptology ePrint Archive*
558 (2016).
- 559 [18] P. Jauernig, A.-R. Sadeghi, E. Stapf, Trusted execution environments: Proper-
560 ties, applications, and challenges, *IEEE Secur. Priv.* 18 (2020) 56–60.
- 561 [19] R. L. Rivest, A. Shamir, L. M. Adleman, A method for obtaining digital
562 signatures and public-key cryptosystems (reprint), *Commun. ACM* 26 (1983)
563 96–99.
- 564 [20] C. Cachin, R. Guerraoui, L. E. T. Rodrigues, *Introduction to Reliable and
565 Secure Distributed Programming (2. ed.)*, Springer, 2011.
- 566 [21] H. Birkholz, D. Thaler, M. Richardson, et al., Remote ATtestation procedureS
567 (RATS) Architecture, 2023.
- 568 [22] S. A. Thomas, *SSL and TLS Essentials: Securing the Web*, Wiley, 2000.

- 569 [23] A. J. M. M. Weijters, W. M. P. van der Aalst, A. K. Alves De Medeiros, Process
570 mining with the HeuristicsMiner algorithm, 2006.
- 571 [24] F. Mannhardt, Sepsis cases - event log, 2016. doi:[10.4121/UUID:](https://doi.org/10.4121/UUID:915D2BFB-7E84-49AD-A286-DC35F063A460)
572 [915D2BFB-7E84-49AD-A286-DC35F063A460](https://doi.org/10.4121/UUID:915D2BFB-7E84-49AD-A286-DC35F063A460).
- 573 [25] W. Steeman, BPI challenge 2013, incidents, 2013. doi:[10.4121/UUID:](https://doi.org/10.4121/UUID:500573E6-ACCC-4B0C-9576-AA5468B10CEE)
574 [500573E6-ACCC-4B0C-9576-AA5468B10CEE](https://doi.org/10.4121/UUID:500573E6-ACCC-4B0C-9576-AA5468B10CEE).
- 575 [26] W. M. P. van der Aalst, Verification of workflow nets, in: ICATPN, 1997, pp.
576 407–426.
- 577 [27] J. P. Barrett, The coefficient of determination—some limitations, The American
578 Statistician (1974) 19–20.
- 579 [28] N. Altman, M. Krzywinski, Simple linear regression, Nature Methods (2015)
580 999–1000.
- 581 [29] D. Basile, C. Di Ciccio, V. Goretti, S. Kirrane, Blockchain based resource
582 governance for decentralized web environments, Frontiers in Blockchain
583 (2023) 1141909.
- 584 [30] D. Basile, C. Di Ciccio, V. Goretti, S. Kirrane, A blockchain-driven architecture
585 for usage control in solid, in: ICDCSW, 2023, pp. 19–24.
- 586 [31] C. Di Ciccio, M. Montali, Declarative process specifications: Reasoning,
587 discovery, monitoring, in: Process Mining Handbook, Springer, 2022, pp.
588 108–152.