

# Detection of Batch Activities from Event Logs

Niels Martin<sup>a,b,c,\*</sup>, Luise Pufahl<sup>d,e,\*</sup>, Felix Mannhardt<sup>f,g</sup>

<sup>a</sup>*Research Foundation Flanders (FWO), Egmontstraat 5, 1000 Brussel, Belgium*

<sup>b</sup>*Hasselt University, Research group Business Informatics, Martelarenlaan 42, 3500 Hasselt, Belgium*

<sup>c</sup>*Vrije Universiteit Brussel, Data Analytics Laboratory, Pleinlaan 2, 1050 Brussel, Belgium*

<sup>d</sup>*HPI, University of Potsdam, Prof.-Dr.-Helmert-Str. 2-3, 14482 Potsdam, Germany*

<sup>e</sup>*SBE, Technische Universitaet Berlin, Einsteinufer 17, 10587 Berlin, Germany*

<sup>f</sup>*SINTEF Digital, Department of Technology Management, Postboks 4760 Torgarden, 7465 Trondheim, Norway*

<sup>g</sup>*NTNU, Department of Computer Science, 7491 Trondheim, Norway*

---

## Abstract

Organizations carry out a variety of business processes in order to serve their clients. Usually supported by information technology and systems, process execution data is logged in an event log. Process mining uses this event log to discover the process' control-flow, its performance, information about the resources, etc. A common assumption is that the cases are executed independently of each other. However, batch work – the collective execution of cases for specific activities – is a common phenomenon in operational processes to save costs or time. Existing research has mainly focused on discovering individual batch tasks. However, beyond this narrow setting, batch processing may consist of the execution of several linked tasks. In this work, we present a novel algorithm which can also detect parallel, sequential and concurrent batching over several connected tasks, i.e., subprocesses. The proposed algorithm is evaluated on synthetic logs generated by a business process simulator, as well as on a real-world log obtained from a hospital's digital whiteboard system. The evaluation shows that batch processing at the subprocess level can be reliably detected.

---

\*Corresponding authors

*Email addresses:* [niels.martin@uhasselt.be](mailto:niels.martin@uhasselt.be) (Niels Martin),  
[luise.pufahl@hpi.de](mailto:luise.pufahl@hpi.de) (Luise Pufahl), [felix.mannhardt@sintef.no](mailto:felix.mannhardt@sintef.no) (Felix Mannhardt)

*Keywords:* Business Process, Batch Activity, Batch Processing, Discovery, Process Mining, Batch Mining

---

## 1. Introduction

In providing products or services to customers or clients, organizations carry out a variety of business processes. A business process consists of a set of activities which are executed in an organizational and technical environment to reach a certain business goal [1]. The enactment of business processes is usually supported by a so-called process-aware information system (PAIS) (e.g. an Enterprise Resource Planning System or a Hospital Information System), where process execution data is registered in an event log [2]. Process mining as a research discipline deploys this type of data to provide insights into the nature of a business process to, e.g., discover a process model.

A common assumption in process mining is that cases, also called process instances, are executed independently from each other. However, batch processing is commonly applied to reduce cost or save time by collectively handling several process instances in a group at specific activities in a business process [3]. Especially in operations research, batching of products or customers is well discussed [4, 5]: in an effort to reduce the set-up or execution cost of physical resources, specific types of products are scheduled in batches on a machine, or customers at a service desk are combined in batches before starting a service. Also, process participants, who get assigned tasks for execution and can organize their work on their own, tend to handle their work in batches [6]. For instance, usually a set of invoices is collected first before they are checked instead of approving each one individually.

Although event logs are a valuable information source to identify batch behavior, the detection of so-called *batch activities* in event logs was discussed by only few works, such as [6, 7, 8, 9]. It enables companies to identify existing batching behavior and identify opportunities to integrate batch processing in a more structured way. Batch processing can occur as *parallel* batching (also referred to as *simultaneous* batching) where several items are processed simultaneously (e.g., blood tubes being tested by a laboratory machine), or *sequential* batching where the cases are still processed individually but, by batching a set, setup time/cost can be reduced (e.g., processing a set of invoices one after the other). Martin et al. [6] additionally consider *concurrent* batch work, which can be observed when a resource starts working on a task

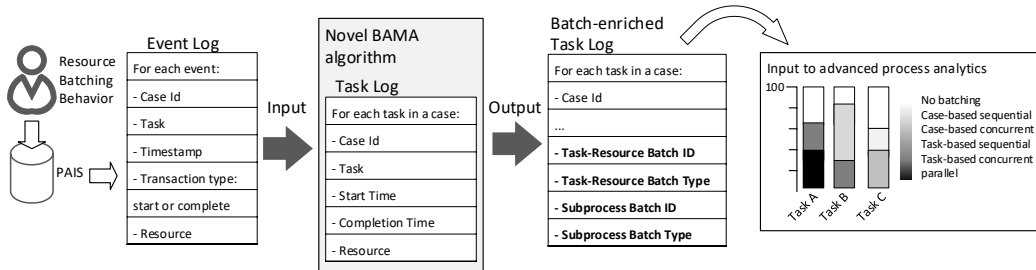


Figure 1: Input and result of the novel BAMA algorithm.

34 and, while executing it, also starts another task instance in parallel. So far,  
 35 the focus of literature on batch mining in the process mining field was situ-  
 36 ated on discovering single batch tasks only, providing only a limited view on  
 37 batch work. Pufahl and Weske [3] observed that batch processing operations  
 38 can span not only over one task, but also over several connected ones.

39 In this paper, we present the novel Batch Activity Mining Algorithm  
 40 (BAMA) which expands existing research by detecting different types of  
 41 batch activities. It extends the algorithm proposed in Martin et al. [6],  
 42 where batching behavior is identified at the level of individual tasks, to the  
 43 more general notion of a batch activity. In BAMA, a batch activity can be  
 44 either non-decomposable (i.e. a task) or can have internal behavior (i.e. a  
 45 subprocess), and for both of them *parallel*, *sequential*, or *concurrent* behavior  
 46 could be detected.

47 Thereby, we also consider that different types of batch behavior can be  
 48 applied during a subprocess. In case of subprocesses, *sequential* batching can  
 49 occur in two versions [10]:

- 50 • with a *task-based* orientation (i.e., cases are executed for each task in  
 51 a subprocess in a sequential fashion) or,
- 52 • with a *case-based* orientation (i.e., first a case needs to finish all tasks  
 53 of a subprocess before the next case of a batch can be started).

54 A similar distinction is made for *concurrent* batching. Human resources can  
 55 apply these types of batch behavior for a set of cases to ease their work.  
 56 Each staff member might apply his/her own strategy for specific tasks. The  
 57 concrete batch execution of activities is often recorded in a PAIS, showing  
 58 the potential of the proposed batch mining algorithm.

59 As shown in Figure 1, the novel BAMA takes as input – as other process  
 60 mining techniques – an event log deduced from a PAIS. Each row represents

61 an event expressing, for instance, the start or completion of a task. BAMA  
62 transforms the event log into a task log, where one row contains all times-  
63 tamps for one task execution. The task log is used for batch mining purposes.  
64 The output of BAMA is a *batch-enriched task log*, a task log to which batch-  
65 ing information is added. For each task executed for a specific case, i.e. for  
66 each task instance, it highlights whether it belongs to a batch by providing  
67 the batch ID and the type of batch behavior. The algorithm provides this at  
68 a task level (i.e., *task-resource*) and at a subprocess level. By its ability to  
69 detect both batch tasks and batch subprocesses, BAMA provides organiza-  
70 tions with rich data-driven insights in batching behavior. This can be used  
71 for advanced analytics, e.g., the occurrences of different batching behavior  
72 for specific tasks can be analyzed. In this paper, the algorithm is formally  
73 described, based on which a prototypical implementation is developed. Be-  
74 sides an evaluation of the algorithm’s effectiveness using synthetic data, this  
75 work also uses real-life data to illustrate how the algorithm’s output can help  
76 practitioners to learn more about batching in their business process.

77 The remainder of this paper is structured as follows. Section 2 presents  
78 the related work. After providing the preliminaries on batch processing in  
79 Section 3, BAMA is introduced in Section 4. In Section 5, the algorithm’s  
80 effectiveness is evaluated using artificial logs. Moreover, it is applied to a  
81 real-world event log obtained from a hospital’s digital whiteboard system.  
82 The paper ends with a discussion in Section 6 and a conclusion in Section 7.

## 83 2. Related Work

84 Batch processing has been studied in various research fields. This section  
85 provides an overview of prior work related to (i) batch processing in oper-  
86 ations management/research, (ii) batch activity modeling, and (iii) batch  
87 mining.

88 *Batch processing in operations management/research.* Batch processing is a  
89 highly relevant topic in operations management/research to save costs or  
90 time by processing or transporting identical jobs in batches, or by utilizing  
91 the same machine/tool set-up [11]. In the following, we do not aim to provide  
92 an exhaustive literature study, but highlight some of the research topics.

93 The order batching problem, which focuses on grouping orders in batches,  
94 is heavily studied in the warehousing context [12, 13, 14, 15]. To this end,  
95 optimization techniques such as integer programming [14, 15] or tabu search

96 [13] are used. Research on this topic is also conducted in other sectors such  
97 as the transportation [16, 17, 18] and the steel sector [19].

98 The scheduling of jobs in batches is also intensively studied as shown  
99 in review papers, such as [20, 21]. Typically, two types of batching are  
100 distinguished: *serial* batching (i.e. similar jobs are scheduled together to save  
101 setup costs, but they are still executed in sequence) and *parallel* batching (i.e.  
102 batches of jobs which can be executed at the same time).

103 Furthermore, operations research offers methods to study the queue of a  
104 batch service (i.e., a service which can handle a group of customers or jobs  
105 in parallel, such as a rollercoaster) to determine, for instance, the expected  
106 waiting time or the expected service time [5]. Thereby, different optimization  
107 policies are proposed. Consider, for instance, the *threshold rule*, which states  
108 that a batch is started when the length of the queue is equal or greater  
109 than a given threshold and the server is free [22]. Several studies (e.g.,  
110 [23]) investigate how to determine the optimal threshold, and the resulting  
111 expected waiting time under varying assumptions.

112 *Batch activity modeling.* In the business process management field, several  
113 approaches were recently developed to capture and specify batch work in  
114 business process models [10, 24, 25]: Pufahl et al. [10, 26] suggest the *batch*  
115 *activity* as a new type of activity with a set of configuration parameters to set  
116 up batch processing for specific activities at design time. The configuration  
117 includes an activation rule, a maximum capacity, and the type of processing  
118 - *parallel* vs. *sequential*. The batch activity can be either a simple task or  
119 an activity with internal behavior (i.e. a subprocess). The performance  
120 evaluation of such batch activities is studied in [27]. Natschläger et al. [24]  
121 propose a so-called *compound activity* for which the optimization goal and  
122 the constraints need to be defined to enable batching. Pflug and Rinderle-Ma  
123 [25] focus on sequential batching and provide an algorithm to rearrange cases  
124 into similar batches which can be processed faster by resources as a way to  
125 organize work for activities with long queues.

126 *Batch mining.* To have more insights in unknown batching behavior in busi-  
127 ness processes, process mining techniques were developed [6, 7, 8, 9, 28].  
128 These techniques use event logs and try to detect tasks for which resources  
129 perform their work in batches. Whereas Wen et al. [7] assume that process  
130 execution data is recorded at the level of a batch, Liu et al. [9], Nakatumba  
131 [8] and Martin et al. [6] use the common event log format in which events

132 are recorded at the level of individual cases. Liu et al. [9] use an event log  
133 to mine association rules for a staff member to recommend work items for  
134 batching, which have been executed collectively in the past. Nakatumba [8]  
135 and Martin et al. [6] propose algorithms to retrieve batching behavior from  
136 an event log. Nakatumba [8] discovers a batch when a time gap of more than  
137 one hour is present between groups of activity instances. This work was  
138 extended by Martin et al. [6] by providing an algorithm to detect different  
139 types of batch work from a log, and providing performance indicators, such  
140 as the size of a batch, waiting time and duration of instances in a batch.

141 Martin et al. [6] not only differentiate between *parallel* and *sequential*  
142 batch processing, but also consider *concurrent* batching as it occurs in prac-  
143 tice that resources work on different instances in a partially overlapping time-  
144 frame, which could be considered as unstructured batch behavior. To distin-  
145 guish between sequential batching and regular queue handling, Martin et al.  
146 [6] assume that all cases need to be present in the queue before the resource  
147 starts processing a batch. When the arrival of a case at a task is recorded,  
148 as is the case in a Q-log in the queue mining field [29], this information can  
149 be immediately retrieved from the data.

150 Different from prior works, Weber et al. [30] provide a pre-processing step  
151 to an existing process discovery algorithm to identify subprocesses in which  
152 multiple instances are initiated in parallel, but executed independently. The  
153 latter holds, for instance, when several reviews are initiated and handled for  
154 one paper submission. Traditionally, these types of subprocesses are called  
155 multi-instance activities [31]. Multi-instance activities allow for running sev-  
156 eral instances of the same activity to handle several sub-items of one process  
157 instance. These instances might be executed independently of each other  
158 and not necessary by the same resource.

159 Besides batch mining and multi-instance initiation, other works, such  
160 as [32, 33], consider a broader spectrum of relations between instances, which  
161 were also called instance-spanning constraints [34]. Whereas Senderovich  
162 et al. [32] use inter-case dependencies to improve the quality of process pre-  
163 dictions with a special focus on case priorities, Winter et al. [33] provide a  
164 technique to detect different kinds of instance-spanning constraints from a  
165 given event log to analyze a process regarding compliance issues. As the fo-  
166 cus is on compliance analysis, different batching behaviors are not analyzed  
167 in depth, which is the goal of our work.

168 So far, the focus of batch mining algorithms was only on detecting single  
169 batch tasks [6, 7, 8, 9] or batching logic at the level of a single task [28].

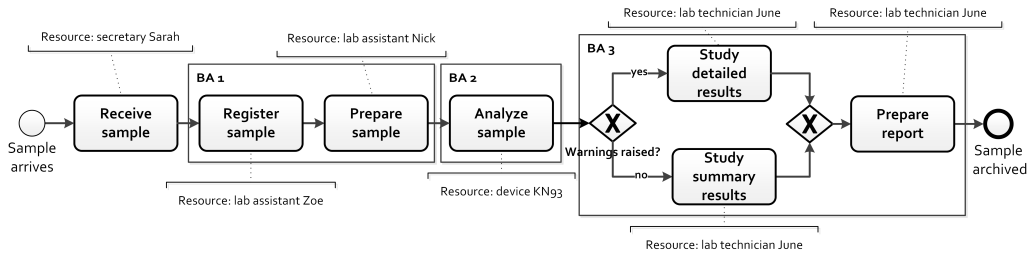


Figure 2: BPMN process diagram of running example 1 - a simplified laboratory process. Using the terminology that will be introduced in Section 3.2, the following batch processing types are included: BA 1 - sequential task-based batching, BA 2 - parallel batching, BA 3 - sequential case-based batching.

170 However, batching is often also done over several linked activities, i.e. in  
 171 subprocesses, which can also have complex behavior, such as parallel activ-  
 172 ities. Given this research gap, this paper presents an algorithm which can  
 173 detect batch behavior at both the task and subprocess level.

### 174 3. Preliminaries

175 The algorithm presented in this paper automatically detects batch activi-  
 176 ties from an event log. Before proceeding to the outline of the algorithm, this  
 177 section outlines some preliminary concepts: Section 3.1 introduces the notion  
 178 of a batch activity, Section 3.2 presents an overview of the batch processing  
 179 types, and Section 3.3 provides an introduction to event logs.

#### 180 3.1. Batch activity

181 Batch processing is an organization of work in which a resource bundles  
 182 cases such that they can be processed as a group [3, 6]. A batch activity is  
 183 a task or a subprocess in which batching behavior is present. We define a  
 184 batch subprocess as a set of connected tasks for which cases are processed in  
 185 batches and all cases stay in the same batch for all subprocess tasks.

186 To illustrate the notion of a batch activity, two running examples are  
 187 introduced. The process models, annotated with the resource responsible to  
 188 perform a task, are visualized as BPMN process diagrams in Figures 2 and  
 189 3.

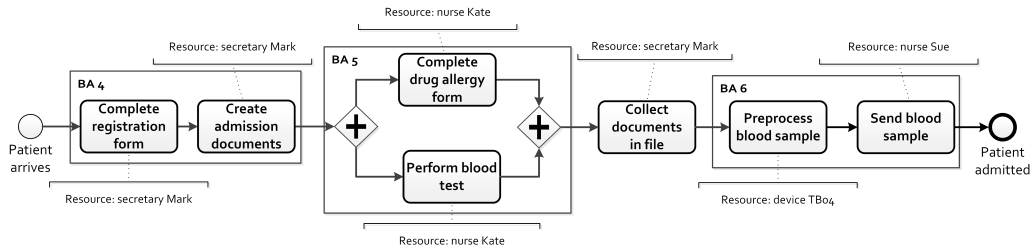


Figure 3: BPMN process diagram of running example 2 - a simplified patient admission process. Using the terminology that will be introduced in Section 3.2, the following batch processing types are included: BA 4 - concurrent task-based batching, BA 5 - concurrent case-based batching, BA 6 - parallel batching.

### 190 3.1.1. Running example 1: A laboratory process

191 The first running example relates to the simplified process that a blood  
 192 sample follows at a hospital laboratory (Figure 2). When a blood sample is  
 193 received by the secretary, it is registered and prepared for the analysis by a  
 194 lab assistant. Afterwards, the actual analysis is conducted by a laboratory  
 195 device. Depending on whether or not warnings are raised during the analysis,  
 196 a lab technician studies either the detailed results or the summary results.  
 197 After studying the results, a report is prepared.

198 Figure 2 shows that three batch activities are present: one consisting of  
 199 a single task (BA 2) and two batch subprocesses (BA 1, BA 3). Besides  
 200 differences in the number of tasks involved, the batch activities also vary  
 201 in terms of the type of batching behavior. In the batch activity consisting  
 202 of task ‘Analyze sample’ (BA 2), multiple samples can be processed by the  
 203 laboratory device simultaneously. In contrast, in BA 1, the lab assistant  
 204 waits until a number of samples is waiting to be processed as this saves setup  
 205 time. Once several samples are available, the lab assistant registers them one  
 206 after the other, after which the same group of samples are prepared one at  
 207 the time. In last batch activity, BA 3, the lab technician will also handle  
 208 samples in groups. She studies the (detailed or summarized) results and  
 209 prepares a report for a particular sample and then immediately proceeds to  
 210 the next sample.

### 211 3.1.2. Running example 2: A patient admission process

212 The second running example focuses on a simplified patient admission  
 213 process at a hospital unit (Figure 3). When a patient arrives, a registra-  
 214 tion form needs to be completed, after which the admission documents are



215 created. Afterwards, a drug allergy form is completed and a blood test is  
216 performed, followed by the collection of all documents in the patient’s file.  
217 When a number of blood samples have been collected, the samples are pre-  
218 processed using a specialized device and sent to the laboratory for further  
219 examination.

220 As shown in Figure 3, three batch subprocesses (BA 4, BA 5, BA 6) are  
221 present. BA 4 contains tasks ‘*Complete registration form*’ and ‘*Create ad-*  
222 *mission documents*’. While a particular patient is filling out the registration  
223 form, the secretary can already start registering another patient. Afterwards,  
224 while the admission files are being printed, the secretary can already compile  
225 the required documents for the next patient. The next batch subprocess, BA  
226 5, consists of tasks ‘*Complete drug allergy form*’ and ‘*Perform blood test*’.  
227 While the patient is filling out a drug allergy form, the nurse already starts  
228 performing the blood test for this patient. Only after several blood samples  
229 are available, BA 6 is executed. This involves preprocessing multiple blood  
230 samples simultaneously, after which all samples are sent to the laboratory at  
231 once for additional analyses.

232 From these running examples, it follows that several types of batching  
233 behavior can be distinguished, which will be introduced in the next subsec-  
234 tion.

### 235 3.2. Batch processing types

236 BAMA will distinguish five types of batch processing. These differ ac-  
237 cording to two dimensions: (i) the time relationship between batched cases,  
238 and (ii) the task- or case-based orientation of a batch. With respect to the  
239 *time relationship* between the execution of tasks on batched cases, a distinc-  
240 tion is made between parallel, sequential and concurrent batching. This is  
241 consistent with Martin et al. [6], where batching behavior is considered at  
242 the level of individual tasks. However, as our algorithm is more generic and  
243 also considers batch detection at the level of a subprocess, an additional dis-  
244 tinction is made between *task-based batching and case-based batching* [10].  
245 While work is organized around tasks in the former (i.e. the same task is  
246 executed for a group of cases before switching to another task), the latter  
247 centers around cases (i.e. several tasks are performed for a particular case  
248 before switching to another case). When both dimensions are taken into  
249 account, five types of batching behavior can be distinguished. Figure 4 il-  
250 lustrates these batch processing types considering three tasks (A, B, C) and  
251 three batched cases (c1, c2, c3).

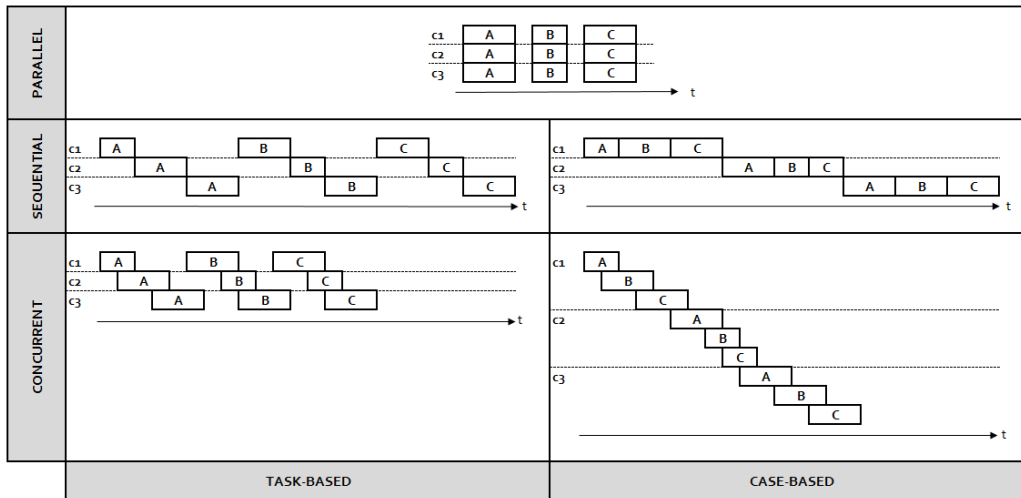


Figure 4: Batch processing types

- 252 • **Parallel batching.** In parallel batching, a set of cases is processed  
 253 at exactly the same time for one or multiple consecutive tasks. For  
 254 each task, all cases are processed by the same resource. In Figure 4, a  
 255 resource first performs task A on cases 1, 2 and 3 simultaneously, and  
 256 does the same for tasks B and C. The task ‘*Analyze sample*’ in Figure 2  
 257 (BA 2) also processes several cases in parallel. In Figure 3, BA 6 is an  
 258 example of two connected activities where parallel batching prevails.
- 259 • **Sequential task-based batching.** In sequential batching, a resource  
 260 handles a group of cases (almost) immediately after each other to save  
 261 set-up time/costs (e.g., the time required to get familiar with a type  
 262 of task). In an effort to differentiate sequential batching from regular  
 263 queue handling, a set of cases can only form a sequential batch when  
 264 all of them are available for the resource before it starts processing  
 265 the batch. In sequential task-based batching, a task will be performed  
 266 on a set of cases sequentially, after which other tasks will sequentially  
 267 be performed to that same set of cases. As with parallel batching, all  
 268 cases are processed by the same resource for a particular task. Figure 4  
 269 shows that task A is performed on the three cases sequentially, which  
 270 is immediately followed by task B and, finally, task C. Batch activity  
 271 BA 1 in Figure 2 is an example of sequential task-based batching.

- 272 • **Sequential case-based batching.** In sequential case-based batching,  
273 a resource will perform a series of tasks sequentially for a particular  
274 case, and will (almost) immediately perform this same task sequence  
275 for other cases. In Figure 4, it is illustrated that tasks A, B and C are  
276 first performed for case 1 first, followed by case 2 and case 3. Batch  
277 activity BA 3 in Figure 2 exemplifies sequential case-based batching  
278 behavior.
  
- 279 • **Concurrent task-based batching.** When concurrent batching pre-  
280 vails, there is a partial overlap in time between the execution of tasks  
281 on a case. This implies that a resource can start executing a new task  
282 on a case before the current task is finished, or that a new case can  
283 be started while the current one has not fully been processed. In con-  
284 current task-based batching, a task is performed for all batched cases,  
285 followed by the other tasks, with the characterizing partial time over-  
286 laps being present. For each task, all cases are processed by the same  
287 resource. Concurrent task-based batching is illustrated in Figure 4,  
288 where work is organized around tasks A, B and C. In Figure 3, BA  
289 4 provides an example of concurrent task-based batching: while the  
290 patient is filling out the registration form, the secretary can already  
291 start the registration of another patient. Similarly, while the admission  
292 documents for one patient are being printed, the secretary can already  
293 start creating the admission documents for another patient.
  
- 294 • **Concurrent case-based batching.** When concurrent case-based  
295 batching prevails, the batch is organized around cases, implying that  
296 the tasks in the batch activity are performed on a particular case be-  
297 fore proceeding to the next case. As shown in Figure 4, the partial  
298 overlap in time between task instances is present both within a case  
299 and between cases. Concurrent case-based batching is also illustrated  
300 by BA 5 in Figure 3 because a nurse can start performing a blood test  
301 for a patient while this patient is completing a drug allergy form.

302 Four remarks need to be made with respect to the batch processing types  
303 outlined above. Firstly, no task-based and case-based variant of parallel  
304 batching is specified. Parallel batching implies that a resource can perform  
305 a specific task on a set of cases simultaneously. From this, it follows that  
306 it is, by definition, task-oriented. When parallel batching would imply that

307 one resource performs several tasks during exactly the same timeframe, these  
308 tasks would, in fact, constitute a single task.

309 Secondly, the requirements regarding resource involvement differ depend-  
310 ing on the batch processing type under consideration. For parallel batching  
311 and sequential/concurrent task-based batching, all instances of a particular  
312 task have to be executed by the same resource. The requirements are more  
313 strict for sequential/concurrent case-based batching as all batched task in-  
314 stances need to be executed by the same resource. For all batch processing  
315 types, it holds the resource cannot be linked to instances which are not part  
316 of a batch while the resource is processing that batch.

317 Thirdly, for case-based batching, the time relation between task instances  
318 associated to a particular case is dominant to distinguish between sequential  
319 and concurrent case-based batching. Consider, for example that a concurrent  
320 relationship holds for task instances for a particular case, but the cases in the  
321 batch are handled sequentially. Under these circumstances, the batch will be  
322 considered as a concurrent case-based batch.

323 Finally, in batch subprocesses, the parallel, sequential or concurrent rela-  
324 tionship between batched instances might differ across tasks (for task-based  
325 batching) or across cases (for case-based batching). For instance, a task of  
326 a task-based batch subprocess can be executed in parallel and another task  
327 in a sequential way. Similarly, the task instances included in a case-based  
328 batch subprocess might be executed sequentially for one case and concur-  
329 rently for another case. In such situations, the detected batch subprocess  
330 will be labeled as a hybrid task-based or case-based batch subprocess. Hy-  
331 brid batching is not included as an autonomous batch processing type as  
332 it merely consists of a combination of different types of batching behavior  
333 within a batch subprocess.

### 334 3.3. *Event log*

335 An event log is the data source that BAMA will use. It contains process  
336 execution information originating from a PAIS [2]. An event log is composed  
337 of a series of events reflecting, e.g., the start or completion of an activity.  
338 An excerpt of the event log from the laboratory process shown Figure 2 is  
339 represented in Table 1. The first line indicates that lab assistant Zoe started  
340 the task ‘*Register sample*’ for sample 9845 on January, 14th at 11:22:33. She  
341 completed this task at 11:26:04, as shown in the second line.

342 The example provided in Table 1 outlines the minimal event log require-  
343 ments to apply the algorithm presented in this paper. The event log should

Table 1: Illustration of event log structure (running example 1)

case id	timestamp	task	transaction type	resource
...	...	...	...	...
9845	14/01/2019 11:22:33	Register sample	start	Lab assistant Zoe
9845	14/01/2019 11:26:04	Register sample	complete	Lab assistant Zoe
9852	14/01/2019 11:26:04	Register sample	start	Lab assistant Zoe
9852	14/01/2019 11:30:21	Register sample	complete	Lab assistant Zoe
9845	14/01/2019 11:30:21	Prepare sample	start	Lab assistant Nick
9893	14/01/2019 11:36:17	Receive sample	start	Secretary Sarah
9845	14/01/2019 11:37:58	Prepare sample	complete	Lab assistant Nick
9852	14/01/2019 11:37:58	Prepare sample	start	Lab assistant Nick
9893	14/01/2019 11:38:12	Receive sample	complete	Secretary Sarah
9852	14/01/2019 11:46:11	Prepare sample	complete	Lab assistant Nick
...	...	...	...	...

344 be an ordered set of events related to a particular case and task. Moreover,  
 345 the timestamp, resource and transaction type needs to be recorded for each  
 346 event. Two transaction types are assumed to be registered: the start and  
 347 the completion of a task, both transitions of the XES-lifecycle extension [35].  
 348 Each start event should have a corresponding complete event with the same  
 349 resource being associated to both events.

350 Organizations usually have multiple resources being able to execute the  
 351 same task. Hence, resource information is needed to identify which resource  
 352 has worked on a group of of cases for this task. If resource information is  
 353 absent and multiple resources perform the same task for different cases, this  
 354 is parallel processing of cases and not batching behavior. The start and end  
 355 time of a task are relevant for identifying the types of batch behavior. With  
 356 only one timestamp for a task, parallel, sequential and concurrent behavior  
 357 are not distinguishable.

358 The aforementioned event log requirements will be formalized below.  
 359 Throughout the paper, we generally use upper case letters for sets and lower  
 360 case letters for elements of sets. So, if  $C$  denotes the set of case identifiers,  
 361 we use  $c \in C$  for a specific case identifier. Table 2 summarizes the key  
 362 mathematical notation which will be used in the paper.

363 **Definition 1** (Event log requirements). *Let  $C$  be a set of case identifiers.*  
 364 *Let  $L$  be a set of task labels. Let  $R$  be a set of resource identifiers. An event*  
 365 *log  $E$  contains a set of events  $e$ . Each event  $e$  is represented as tuple of*  
 366 *attributes  $e = (c, t, r, \tau, \varphi)$  with:*

- 367 •  $c \in C$  represents the case identifier,

Table 2: Summary of mathematical notation

Symbol	Description
$E/e$	Set of events / event
$C/c$	Set of case identifiers / case
$L/l$	Set (element) of task labels / task
$R/r$	Set (element) of resource identifiers / resource
$\mathbb{R}^+$	Set of positive real numbers
$\tau$	Timestamp attribute of an event
$\varphi$	Transaction type
$\#_a(e)$	Function returning the value of attribute $a$ of event $e$
$E_{start}$	Set of events with transaction type <i>start</i>
$E_{complete}$	Set of events with transaction type <i>complete</i>
$T$	Set of task instance (Task log)
$T_{ex}$	Set of extended task instances (Extended task log)
$\tau_{arrival}$	Timestamp attribute of the task instance's arrival
$\tau_{start}$	Timestamp attribute of the task instance's start
$\tau_{complete}$	Timestamp attribute of the task instance's completion
$\chi$	Set of task labels
$T_\chi$	Subset of an extended task log having task labels $\chi$
$\phi$	Set of case identifiers
$T_\phi$	Subset of an extended task log having case identifiers $\phi$
$T_{be}$	Set of batch-enriched task instances (Batch-enriched task log)

- 368 •  $l \in L$  is the task label,
- 369 •  $r \in R$  represents the resource identifier,
- 370 •  $\tau \in \mathbb{R}^+$  reflects the timestamp, and
- 371 •  $\varphi \in \{start, complete\}$  refers to the transaction type.

372 We use the notation  $\#_a(e)$  to access the value for attribute  $a$  for event  $e$ .

373 Moreover, every start event should have an accompanying complete event.

374 In formal terms, let  $E_{start} = \{e \in E \mid \#_\varphi = start\}$  and  $E_{complete} = \{e \in$   
375  $E \mid \#_\varphi = complete\}$  be a partition of event log  $E$ . Then, there exists a  
376 bijective function  $\omega : E_{start} \rightarrow E_{complete}$ , such that  $\forall e \in E_{start} : \#_c(e) =$   
377  $\#_c(\omega(e)) \wedge \#_l(e) = \#_l(\omega(e)) \wedge \#_r(e) = \#_r(\omega(e)) \wedge \#_\tau(e) \leq \#_\tau(\omega(e))$ .

#### 378 4. Batch Activity Mining Algorithm

379 This section outlines the Batch Activity Mining Algorithm (BAMA),  
380 which automatically identifies batching behavior in an event log, both at  
381 the task and the subprocess level. Section 4.1 provides a general overview  
382 of the algorithm, after which its key steps are outlined in more detail in  
383 Sections 4.2- 4.6. The prototypical implementation is briefly discussed in

384 Section 4.7. Section 4.8 provides some pointers on how the output of BAMA  
385 can be used to gain a rich understanding in batching behavior.

#### 386 4.1. General overview

387 BAMA aims to identify batching behavior at two distinct levels: (i) at  
388 the level of individual tasks and (ii) at the level of subprocesses. As shown  
389 in Figure 5, an event log is used as an input. This event log is converted to a  
390 task log by mapping start events to their corresponding completion events.  
391 The task log is used for batch identification purposes.

392 To achieve BAMA’s goals, three phases can be distinguished in the al-  
393 gorithm. In the first phase, the task log is extended with batch formation  
394 insights at the task-resource level, i.e. batches consisting of task instances of  
395 a particular task executed by a particular resource. For instance: when sev-  
396 eral blood samples are analyzed by ‘*Device KN93*’ simultaneously, a parallel  
397 batch is formed at the task-resource level for task-resource combination ‘*An-*  
398 *alyze sample*’ - ‘*Device KN93*’. For this phase of the algorithm, the algorithm  
399 presented in Martin et al. [6] is used.

400 While the first phase of the algorithm focuses on the identification of  
401 batching behavior at the level of individual tasks, the other two phases tar-  
402 get the discovery of batch subprocesses. To this end, the extended task  
403 log is split into two subsets using the batch insights at the task-resource  
404 level: subset 1 containing candidates for a parallel or task-based sequen-  
405 tial/concurrent batch subprocess, and subset 2 with candidates for a case-  
406 based sequential/concurrent batch subprocess. The former is used as input  
407 for batch subprocess discovery in the second phase of the algorithm and the  
408 latter in the third phase.

409 Merging the outcomes of the second and third phase generates a batch-  
410 enriched task log. Compared to the original task log, a batch-enriched task  
411 log contains four additional columns providing batching information: two  
412 columns highlighting whether the task instance is part of a batch at the  
413 task-resource level (one column carrying a batch identifier and another one  
414 containing the batch type), and two columns reflecting its membership of a  
415 batch subprocess (again, a batch identifier column and batch type column are  
416 added). These columns enable analysts to retrieve rich insights in batching  
417 behavior starting from real-life process execution data.

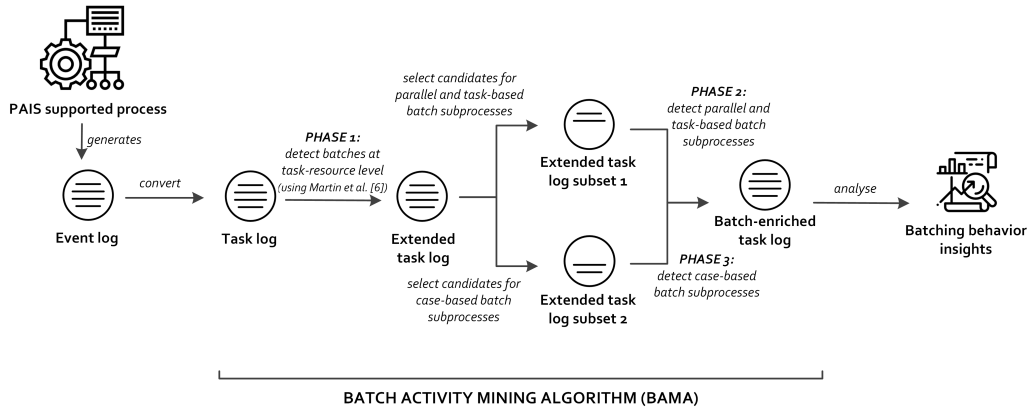


Figure 5: Overview of the Batch Activity Mining Algorithm (BAMA)

418 *4.2. Convert to task log*

419 The input of the algorithm is an event log, of which the minimal require-  
 420 ments were outlined in Section 3.3. As the detection of batching behavior  
 421 requires studying the time relationship between task instances, the event log  
 422 needs to be converted to a task log. In a task log, each row represents a task  
 423 instance, i.e. the execution of a particular task by a particular resource in a  
 424 particular case.

425 To convert the event log to a task log, each start event is mapped to  
 426 its corresponding completion event, where the latter refers to the complete  
 427 event which is linked to the same case, activity and resource. For instance:  
 428 the first and second line in Table 1 will occur as a single line in the task log.  
 429 When multiple start and complete events are recorded for a particular task-  
 430 resource-case combination, the first occurring start event will iteratively be  
 431 mapped to the first occurring unmapped completion event. For more complex  
 432 mappings, techniques such as the one described by Baier et al. [36] can be  
 433 used. The task log created from the event log excerpt in Table 1 is shown in  
 434 Table 3.

435 **Definition 2** (Task log). *Let  $E$  represent an event log and let  $T$  be a task*  
 436 *log. To convert  $E$  into  $T$ , a mapping function  $m_1$  is applied. Function  $m_1 :$*   
 437  *$E \rightarrow T$  maps the events in  $E$  to task instances in task log  $T$ . This is achieved*  
 438 *by mapping start events  $e_s \in E_{start}$  to their corresponding complete events*  
 439  *$\omega(e_s) \in E_{complete}$  to combine start and completion times. Consequently,  $T$*   
 440 *consists of a set of task instances  $t$ . Each task instance  $t$  is represented as*



Table 3: Illustration of task log structure (running example 1)

case id	task	resource	$\tau_{start}$	$\tau_{complete}$
...	...	...	...	...
9845	Register sample	Lab assistant Zoe	14/01/2019 11:22:33	14/01/2019 11:26:04
9852	Register sample	Lab assistant Zoe	14/01/2019 11:26:04	14/01/2019 11:30:21
9845	Prepare sample	Lab assistant Nick	14/01/2019 11:30:21	14/01/2019 11:37:58
9893	Receive sample	Secretary Sarah	14/01/2019 11:36:17	14/01/2019 11:38:12
9852	Prepare sample	Lab assistant Nick	14/01/2019 11:37:58	14/01/2019 11:46:11
...	...	...	...	...

441 a tuple of attributes  $t = (c, l, r, \tau_{start}, \tau_{complete})$ , where  $\tau_{start}$  refers to the start  
 442 time, and  $\tau_{complete}$  is the completion time.

#### 443 4.3. Phase 1: detect batches at the task-resource level

444 The algorithm’s first phase focuses on the detection of batching behavior  
 445 at the task-resource level, i.e. on identifying groups of task instances which  
 446 are processed as a batch at a particular task by a particular resource. This  
 447 task-resource level is purposefully selected as it conveys valuable insights in  
 448 the way in which a specific resource executes a task.

449 Besides the useful information it offers in itself, batching behavior at the  
 450 task-resource level is also a prerequisite for the presence of some types of  
 451 batch subprocesses. From the outline of the batch processing types in Sec-  
 452 tion 3.2, it follows that parallel and task-based sequential/concurrent batch  
 453 subprocesses require that batching behavior is detected at the task-resource  
 454 level. To this end, the first phase of BAMA involves detecting parallel (par),  
 455 sequential (seq) or concurrent (conc) batches at the task-resource level.

456 To detect batches at the task-resource level, the algorithm presented in  
 457 Martin et al. [6] is used. This algorithm will group task instances for which a  
 458 batching relationship holds by marking these instances with the same batch  
 459 number. The parallel, sequential or concurrent character of the batch is also  
 460 added.

461 As follows from the description of sequential batching in Section 3.2, the  
 462 arrival time of a case at a task is used to distinguish between sequential  
 463 batching and regular queue handling. Consequently, each task instance in  
 464 the task log needs to be enriched with the arrival time (i.e. the time at which  
 465 a task instance is enabled). The arrival time differs from the start time of a  
 466 task instance when queues are formed because of limited resources to perform  
 467 a particular task. In a Q-log, which is an event log containing queue-related

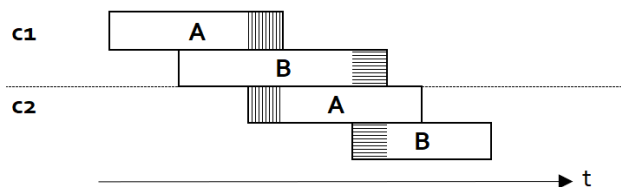


Figure 6: Illustration of concurrent batching at the task-resource level

468 events such as queue arrival events [29], arrival times are explicitly recorded  
 469 in the event log. Hence, they can just be included when the event log is  
 470 converted to a task log and no further efforts are required. When, instead,  
 471 arrival times are unknown, they can be imputed using a suitable heuristic.  
 472 For example, the task arrival time of a case can be approximated by the end  
 473 of its preceding task. This requires insights in the process control-flow, which  
 474 can be gathered from the event log using an existing control-flow discovery  
 475 algorithm. As control-flow discovery is beyond the scope of this paper, the  
 476 reader is referred to, e.g., van der Aalst [2], De Weerd et al. [37], and Augusto  
 477 et al. [38] for more background on this topic.

478 When batches are detected at the task-resource level, these task instances  
 479 are candidates to be part of a parallel or task-based sequential/concurrent  
 480 batch subprocess (cf. phase 2 of the algorithm). Conversely, task instances  
 481 which are not batched at the task-resource level are candidates to be part  
 482 of a case-based sequential/concurrent batch subprocess (cf. phase 3 of the  
 483 algorithm). The latter follows from the observation that case-based batching  
 484 does not require batching behavior at the task-resource level, as shown at  
 485 the right side of Figure 4. However, detected concurrent and sequential  
 486 batches at the task-resource level are also included as candidates for case-  
 487 based batch subprocesses. To illustrate why this is the case, consider the  
 488 illustrations of a concurrent case-based batch subprocess in Figure 6 and  
 489 Figure 7 (where all instances are performed by the same resource). Due to  
 490 the strong time overlap between the instances, the two instances of task A  
 491 form a concurrent batch at the task-resource level in Figure 6 and a sequential  
 492 batch in Figure 7. A similar remark holds for the two instances of task B.  
 493 These simple illustrations show why concurrent and sequential batches at  
 494 the task-resource level are also candidates to be part of a case-based batch  
 495 subprocess.

496 From what said, it follows that the first phase of the algorithm involves  
 497 extending the task log with three columns: the task arrival time, a task-

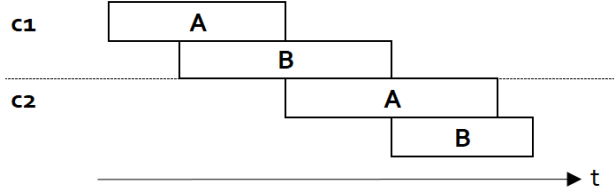


Figure 7: Illustration of sequential batching at the task-resource level

498 resource batch identifier and a task-resource batch type. When extending  
 499 the task log in Table 3, the log in Table 4 is obtained. Table 4 shows that  
 500 two sequential batches are detected at the task-resource level: batch 138  
 501 containing two instances of task ‘Register sample’ and batch 374 consisting  
 502 of two instances of the task ‘Prepare sample’.

503 **Definition 3** (Extended task log). *Let  $T$  represent a task log and let  $T_{ex}$  be*  
 504 *an extended task log. To convert  $T$  into  $T_{ex}$ , two functions  $m_2$  and  $m_3$  are*  
 505 *sequentially applied:*

- 506 •  $m_2 : T \rightarrow T'$  adds the arrival time  $\tau_{arrival}$  for each task instance. In  
 507 case the arrival times are recorded in the event log,  $m_2$  can be merged  
 508 with  $m_1$  in Definition 2.
- 509 •  $m_3 : T' \rightarrow T_{ex}$  assigns a task-resource batch identifier  $b_{tr,id}$  and a task-  
 510 resource batch type  $b_{tr,t}$  to all task instances  $t \in T'$  with  $b_{tr,id} \in \mathbb{N}$  and  
 511  $b_{tr,t} \in \{par, seq, conc\}$ . The values of  $b_{tr,id}$  and  $b_{tr,t}$  are shared among  
 512 the task instances  $t$  which are part of the same batch at the task-resource  
 513 level <sup>1</sup>.

514 After the application of  $m_2$  and  $m_3$ , the extended task log  $T_{ex}$  consists of a  
 515 set of batched task instances  $t_b$ . Each batched task instance is represented as  
 516 a tuple of attributes  $t_b = (c, l, r, \tau_{arrival}, \tau_{start}, \tau_{complete}, b_{tr,id}, b_{tr,t})$ .

#### 517 4.4. Phase 2: detect parallel and task-based sequential/concurrent batch sub- 518 processes

519 Using the task instances which are batched at the task-resource level as an  
 520 input (cf. subset 1 in Figure 5), the second phase involves detecting parallel  
 521 and task-based sequential/concurrent batch subprocesses.

---

<sup>1</sup>The application of function  $m_3$  involves the use of the method defined by Martin et al. [6] to identify batches at the task-resource level.

Table 4: Illustration of extended task log structure (running example 1)

case id	task	resource	$\tau_{arrival}$	$\tau_{start}$	$\tau_{complete}$	task-resource batch id	task-resource batch type
...	...	...	...	...	...	...	...
9845	Register sample	Lab assistant Zoe	14/01/2019 10:17:38	14/01/2019 11:22:33	14/01/2019 11:26:04	138	seq
9852	Register sample	Lab assistant Zoe	14/01/2019 10:32:44	14/01/2019 11:26:04	14/01/2019 11:30:21	138	seq
9845	Prepare sample	Lab assistant Nick	14/01/2019 11:26:04	14/01/2019 11:30:21	14/01/2019 11:37:58	374	seq
9893	Receive sample	Secretary Sarah	14/01/2019 11:30:21	14/01/2019 11:36:17	14/01/2019 11:38:12	-	-
9852	Prepare sample	Lab assistant Nick	14/01/2019 11:14:17	14/01/2019 11:37:58	14/01/2019 11:46:11	374	seq
...	...	...	...	...	...	...	...

522 To identify these subprocesses, two conditions need to be verified. Firstly,  
523 it is checked whether there are different batches at the task-resource level  
524 having exactly the same composition in terms of cases. For instance, in  
525 Table 4, the detected batches for ‘*Register sample*’ - ‘*Lab assistant Zoe*’ and  
526 ‘*Prepare sample*’ - ‘*Lab assistant Nick*’ contain the same cases 9845 and  
527 9852. Secondly, when such batches are present, it is verified whether, for  
528 each case, these tasks immediately follow each other (i.e. that no other tasks  
529 have been executed in between). When both conditions are fulfilled, a batch  
530 subprocess has been detected. All task instances included in the subprocess  
531 will be marked with the same batch subprocess identifier.

532 To determine whether the detected batch subprocess is a parallel or task-  
533 based sequential/concurrent batch subprocess, the batch types at the task-  
534 resource level are taken into consideration. When all tasks included in a  
535 batch subprocess are executed as a parallel batch, the batch subprocess type  
536 will be parallel. Similarly, when all included tasks are executed as a sequen-  
537 tial or concurrent batch, the batch subprocess type will be sequential and  
538 concurrent, respectively. In case a combination of batch types is present  
539 at the task-resource level, e.g. one task is executed in parallel and another

540 one sequential, the batch subprocess will be referred to as a hybrid batch  
 541 subprocess.

542 With respect to sequential batching at the task-resource level, Martin  
 543 et al. [6] allow for a time gap between the end of a particular task instance  
 544 and the start of the next one. A non-zero time gap can accommodate, e.g.,  
 545 for some set-up time required between cases in practical applications [6]. In  
 546 the detection of sequential task-based subprocesses, a similar time gap can  
 547 be defined. The time gap expresses the maximal tolerable time period that  
 548 can elapse between the end of the last task instance of a particular task in  
 549 the subprocess and the start of the first instance at the next task. Note that  
 550 the requirement that no other tasks can be executed in this time period still  
 551 holds. An appropriate value for the time gap is context-specific and, hence,  
 552 requires domain knowledge. The event log can support domain experts by  
 553 providing insights in the time gaps prevailing in reality.

554 We now formally define the different batch subprocess types. First, re-  
 555 quirements common for all task-based batch subprocess types are identi-  
 556 fied. Then, we define parallel batching and sequential/concurrent task-based  
 557 batching. To clarify the definitions, an illustration for each task-based batch-  
 558 ing type is provided, based on the running examples introduced earlier. For  
 559 each batching type, both its footprint in the batch-enriched task log and a  
 560 visualization is included.

561 **Definition 4** (Base task-based batch conditions). *Let  $\chi \subseteq L$  represent a set*  
 562 *of task labels. The base batch conditions common to all task-based batch types*  
 563 *that a set of task instances  $T_\chi \subseteq T_{ex}$  needs to fulfill are:*

- 564
1. *all task instances refer to task labels in  $\chi$ , i.e.,  $\forall t \in T_\chi \#l(t) \in \chi$ ;*
  2. *the same number<sup>2</sup> and at least two instances for each task label are recorded, i.e.:*

$$\forall_{l_1, l_2 \in \chi} (|\{t \in T_\chi \mid \#l(t) = l_1\}| = |\{t \in T_\chi \mid \#l(t) = l_2\}| \geq 2);$$

3. *given  $1 \leq k \leq n-1$  and the sequence of task instances  $\langle t_1, \dots, t_k, \dots, t_n \rangle \in T_\chi^*$  with  $\#_{\tau_{start}}(t_k) \leq \#_{\tau_{start}}(t_{k+1})$  and  $t_k \neq t_{k+1}$  all task instances with*

---

<sup>2</sup>Note that we simplify the definition by restricting subprocesses to not repeat individual tasks. However, this limitation can be lifted by a suitable renaming of the task labels such that each repetition is uniquely labeled.

the same task label are started before proceeding to task instances with another task label:

$$\#_l(t_k) \neq \#_l(t_{k+1}) \implies \forall_{k+1 \leq j \leq n} (\#_l(t_j) \neq \#_l(t_k))$$

4. the same set of cases is involved in the batch, i.e.:

$$\begin{aligned} & \forall_{l_1, l_2 \in \chi} (\{c \in C \mid \exists t \in T_\chi : \#_c(t) = c \wedge \#_l(t) = l_1\} \\ & = \{c \in C \mid \exists t \in T_\chi : \#_c(t) = c \wedge \#_l(t) = l_2\}); \end{aligned}$$

- 565 5. the same resource executes tasks with the same task label in the full  
566 batch, i.e.,

$$\forall_{t_a, t_b \in T_\chi} (\#_l(t_a) = \#_l(t_b) \implies \#_r(t_a) = \#_r(t_b)).$$

567 **Definition 5** (Parallel batch). Let  $\chi \subseteq L$  represent a set of task labels. A  
568 parallel batch is composed of a set of task instances  $T_\chi \subseteq T_{ex}$  fulfilling the  
569 base task-based batch conditions (cf. Definition 4) and for all task instances  
570  $t_a, t_b \in T_\chi$  the following additional conditions hold:

1. tasks with the same label are performed in parallel:

$$\begin{aligned} \#_l(t_a) = \#_l(t_b) & \iff \#_{\tau_{start}}(t_a) = \#_{\tau_{start}}(t_b) \\ & \wedge \#_{\tau_{complete}}(t_a) = \#_{\tau_{complete}}(t_b); \end{aligned}$$

2. tasks with different labels are non-overlapping:

$$\#_l(t_a) \neq \#_l(t_b) \iff \#_{\tau_{complete}}(t_a) < \#_{\tau_{start}}(t_b);$$

3. no other task instance performed by the same resource, which is not included in the parallel batch, can be started or completed when processing the batched instances, i.e., given:

$$\begin{aligned} T_r = \{t \in (T_{ex} \setminus T_\chi) \mid & \exists_{t_a, t_b \in T_\chi} (\#_r(t_a) = \#_r(t) = \#_r(t_b) \\ & \wedge (\#_{\tau_{start}}(t_a) \leq \#_{\tau_{start}}(t) \leq \#_{\tau_{complete}}(t_b) \\ & \vee \#_{\tau_{start}}(t_a) \leq \#_{\tau_{complete}}(t) \leq \#_{\tau_{complete}}(t_b))\} \end{aligned}$$

571 we require  $T_r = \emptyset$ .

572 An illustration of a parallel batch subprocess within the context of run-  
573 ning example 2 is provided in Table 5 and visualized in Figure 8. The parallel  
574 batch subprocess consists of cases 9072 and 9080. First, these cases are pro-  
575 cessed in a parallel batch at task ‘Preprocess blood sample’. Afterwards, task  
576 ‘Send blood sample’ is executed in these same cases as a parallel batch. As  
577 ‘Preprocess blood sample’ and ‘Send blood sample’ immediately follow each  
578 other for both cases which indicates that a parallel batch subprocess is de-  
579 tected.

580 **Definition 6** (Sequential task-based batch). *Let  $\chi \subseteq L$  represent a set of*  
581 *task labels. A sequential task-based batch is composed of a set of task instances*  
582  *$T_\chi \subseteq T_{ex}$  fulfilling the base task-based batch conditions (cf. Definition 4) and*  
583 *additionally the following conditions:*

- 584 1. *for all  $t \in T_\chi$  the arrival time of the instance cannot be later than the*  
585 *start time of the first instance for that particular task in the batch, i.e.,*  
586 *for all  $t_x \in T_\chi : (\#_l(t_x) = \#_l(t)) \implies \#_{\tau_{arrival}}(t) \leq \#_{\tau_{start}}(t_x)$ ;*  
587 2. *for  $1 \leq k \leq n-1$  and the sequence<sup>3</sup> of task instances  $\langle t_1, \dots, t_k, \dots, t_n \rangle \in$*   
588  *$T_\chi^*$  with  $\#_{\tau_{start}}(t_k) < \#_{\tau_{start}}(t_{k+1})$  the following conditions cumulatively*  
589 *hold:*

- *task instances referring to the same label are at most separated by a time gap  $\gamma$ , i.e.:*

$$\#_l(t_k) = \#_l(t_{k+1}) \implies (\#_{\tau_{start}}(t_{k+1}) - \#_{\tau_{complete}}(t_k)) \in [0, \gamma]$$

*with  $\gamma \geq 0$ ,*

- *the time gap between the completion of the last instance for a particular task in  $\chi$  and the start of the first instance of the next task in  $\chi$  cannot exceed the tolerated time gap  $\theta$ , i.e.:*

$$\#_l(t_k) \neq \#_l(t_{k+1}) \implies (\#_{\tau_{start}}(t_{k+1}) - \#_{\tau_{complete}}(t_k)) \in [0, \theta]$$

*with  $\theta \geq 0$ .*

---

<sup>3</sup>We denote with  $X^*$  the set of all sequences over set  $X$ . Moreover, we assume in the remainder of this paper that it is possible to obtain a totally ordered sequence of events performed by a single resource. If this is not fulfilled, the event logs can be pre-processed based on a secondary order criteria e.g., the completion time.

Table 5: Illustration of a parallel batch subprocess (running example 2)

case id	task	resource	$\tau_{arrival}$	$\tau_{start}$	$\tau_{complete}$	task-resource batch id	task-resource batch type	batch subprocess id	batch subprocess type
...	...	...	...	...	...	...	...	...	...
9072	Preprocess blood sample	Device TB04	10/01/2019 12:28:02	10/01/2019 13:03:17	10/01/2019 13:07:52	52	par	8	par
9080	Preprocess blood sample	Device TB04	10/01/2019 12:35:02	10/01/2019 13:03:17	10/01/2019 13:07:52	52	par	8	par
...	...	...	...	...	...	...	...	...	...
9072	Send blood sample	Nurse Sue	10/01/2019 11:07:52	10/01/2019 13:13:24	10/01/2019 13:17:04	59	par	8	par
9080	Send blood sample	Nurse Sue	10/01/2019 11:07:52	10/01/2019 13:13:24	10/01/2019 13:17:04	59	par	8	par
...	...	...	...	...	...	...	...	...	...

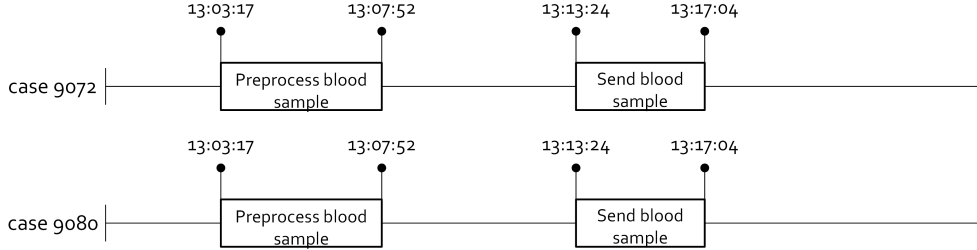


Figure 8: Visualization of a parallel batch subprocess (running example 2)

3. *no other task instance performed by the same resource, which is not included in the sequential task-based batch, can be started or completed when processing the batched instances, i.e., given:*

$$\begin{aligned}
 T_r = \{t \in (T_{ex} \setminus T_\chi) \mid & \exists_{t_a, t_b \in T_\chi} (\#_r(t_a) = \#_r(t) = \#_r(t_b)) \\
 & \wedge (\#_{\tau_{start}}(t_a) \leq \#_{\tau_{start}}(t) \leq \#_{\tau_{complete}}(t_b)) \\
 & \vee \#_{\tau_{start}}(t_a) \leq \#_{\tau_{complete}}(t) \leq \#_{\tau_{complete}}(t_b))\}
 \end{aligned}$$



590 we require  $T_r = \emptyset$ .

591 As illustrated in Table 6 and visualized in Figure 9, a sequential task-  
 592 based batch subprocess is detected for the task instances included in Table 4  
 593 (related to running example 1). More specifically, the batches with task-  
 594 resource batch identifier 138 and 374 fulfill the conditions of a task-based  
 595 sequential batch process as they both consist of the same cases 9845 and  
 596 9852, and the tasks immediately follow each other for both cases. Because  
 597 both tasks are sequential batches at the task-resource level, the task-based  
 598 subprocess is marked as a task-based sequential batch subprocess.

599 **Definition 7** (Concurrent task-based batch). *Let  $\chi \subseteq L$  represent a set  
 600 of task labels. A concurrent task-based batch is composed of a set of task  
 601 instances  $T_\chi \subseteq T_{ex}$  fulfilling the base task-based batch conditions (cf. Defini-  
 602 tion 4) and additionally the following conditions:*

1. given  $1 \leq k \leq n-1$  and the sequence of task instances  $\langle t_1, \dots, t_k, \dots, t_n \rangle \in T_\chi^*$  with  $\#_{\tau_{start}}(t_k) \leq \#_{\tau_{start}}(t_{k+1})$  and  $t_k \neq t_{k+1}$  subsequent task instances are performed concurrently, i.e.:

$$\begin{aligned} \#_{\tau_{start}}(t_k) &\leq \#_{\tau_{start}}(t_{k+1}) < \#_{\tau_{complete}}(t_k) \\ &\wedge (\#_{\tau_{start}}(t_k) \neq \#_{\tau_{start}}(t_{k+1}) \vee \\ &\quad \#_{\tau_{complete}}(t_k) \neq \#_{\tau_{complete}}(t_{k+1})) \end{aligned}$$

2. no other task instance performed by the same resource, which is not included in the concurrent task-based batch, can be started or completed when processing the batched instances, i.e., given:

$$\begin{aligned} T_r = \{t \in (T_{ex} \setminus T_\chi) \mid \exists_{t_a, t_b \in T_\chi} (\#_r(t_a) = \#_r(t) = \#_r(t_b) \\ \wedge (\#_{\tau_{start}}(t_a) \leq \#_{\tau_{start}}(t) \leq \#_{\tau_{complete}}(t_b) \\ \vee \#_{\tau_{start}}(t_a) \leq \#_{\tau_{complete}}(t) \leq \#_{\tau_{complete}}(t_b))\} \end{aligned}$$

603 we require  $T_r = \emptyset$ .

604 An example of a concurrent task-based batch in the patient admission  
 605 process (running example 2) is presented in Table 7 and Figure 10. At  
 606 the task-resource level, the instances of both ‘Complete registration form’  
 607 and ‘Create admission documents’ partially overlap in time. Consequently,  
 608 these instances form concurrent batches at the task-resource level (with task-  
 609 resource batch identifier 36 for ‘Complete registration form’ and 42 for ‘Create

Table 6: Illustration of a sequential task-based batch subprocess (running example 1)

case id	task	resource	$T_{arrival}$	$T_{start}$	$T_{complete}$	task-resource batch id	task-resource batch type	batch subprocess id	batch subprocess type
...	...	...	...	...	...	...	...	...	...
9845	Register sample	Lab assistant Zoe	14/01/2019 10:17:38	14/01/2019 11:22:33	14/01/2019 11:26:04	138	seq	38	seq task-based
9852	Register sample	Lab assistant Zoe	14/01/2019 10:32:44	14/01/2019 11:26:04	14/01/2019 11:30:21	138	seq	38	seq task-based
...	...	...	...	...	...	...	...	...	...
9845	Prepare sample	Lab assistant Nick	14/01/2019 11:26:04	14/01/2019 11:30:23	14/01/2019 11:37:58	374	seq	38	seq task-based
9852	Prepare sample	Lab assistant Nick	14/01/2019 11:14:17	14/01/2019 11:37:58	14/01/2019 11:46:11	374	seq	38	seq task-based
...	...	...	...	...	...	...	...	...	...

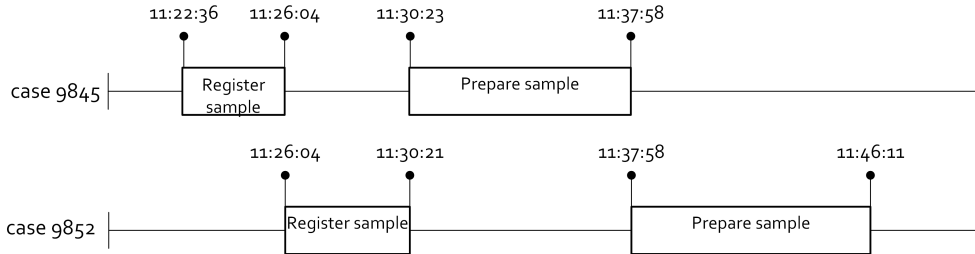


Figure 9: Visualization of a sequential task-based batch subprocess (running example 1)

610 *admission documents*'). As these tasks immediately follow each other for  
 611 both cases, a concurrent task-based batch subprocess is present.

612 In summary, phase 2 of the algorithm can identify sets of task instances  
 613 fulfilling the conditions of either parallel or task-based sequential/concurrent  
 614 batch subprocesses. Two columns are added in this phase: a batch subprocess  
 615 process identifier and a batch subprocess type. In the output, task instances  
 616 belonging to a batch subprocess will share the same batch subprocess identifier.  
 617

Table 7: Illustration of a concurrent task-based batch subprocess (running example 2)

case id	task	resource	$T_{arrival}$	$T_{start}$	$T_{complete}$	task-resource batch id	task-resource batch type	batch subprocess id	batch subprocess type
...	...	...	...	...	...	...	...	...	...
9097	Complete registration form	Secretary Mark	10/01/2019 11:12:18	10/01/2019 11:12:18	10/01/2019 11:19:54	36	conc	4	conc task-based
9098	Complete registration form	Secretary Mark	10/01/2019 11:16:04	10/01/2019 11:16:04	10/01/2019 11:23:31	36	conc	4	conc task-based
9097	Create admission documents	Secretary Mark	10/01/2019 11:19:54	10/01/2019 11:23:31	10/01/2019 11:28:19	42	conc	4	conc task-based
9098	Create admission documents	Secretary Mark	10/01/2019 11:23:31	10/01/2019 11:26:48	10/01/2019 11:32:08	42	conc	4	conc task-based
...	...	...	...	...	...	...	...	...	...

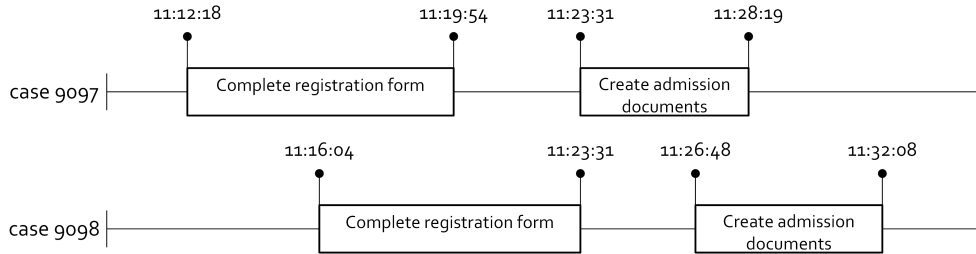


Figure 10: Visualization of a concurrent task-based batch subprocess (running example 2)

618 4.5. Phase 3: detect case-based sequential/concurrent batch subprocesses

619 The algorithm's third phase aims to detect case-based sequential or con-  
 620 current batch subprocesses. To this end, for the reasons outlined earlier, it  
 621 uses the task instances which are not part of a task-resource batch or which  
 622 are included in a concurrent or sequential task-resource batch (cf. subset 2  
 623 in Figure 5). Regarding the inclusion of instances in a concurrent/sequential  
 624 batch at the task-resource level, it should be noted that instances included in

625 a concurrent/sequential task-based batch subprocess are removed from sub-  
626 set 2 after phase 2. This avoids that task instances will be duplicated in the  
627 merged batch-enriched task log after phase 3.

628 To identify case-based batching behavior for a particular set of connected  
629 tasks (i.e. a task subsequence representing a potential subprocess), three  
630 steps are taken. Firstly, *within-case batching* is checked, which focuses on  
631 the presence of a batching relationship between task instances of one partic-  
632 ular case. To illustrate this, consider Table 8, where the subsequence ‘*Study*  
633 *summary results*’ - ‘*Prepare report*’ is considered. As shown in Figure 4,  
634 sequential/concurrent case-based batching requires that a batching relation-  
635 ship is present between the execution of different tasks in a particular case.  
636 This holds for case 9969 in Table 8, where a sequential relationship holds  
637 between the execution of ‘*Study summary results*’ and ‘*Prepare report*’. The  
638 same holds for case 9974.

639 Secondly, when within-case batching is detected, the involved task in-  
640 stances related to a particular case are replaced by a *single aggregated in-*  
641 *stance*. This is a preparatory step for the third step of case-based subprocess  
642 detection. The arrival and start timestamp of this aggregated instance cor-  
643 respond to the earliest arrival and start timestamp of the included instances.  
644 For the complete timestamp, the aggregated instance will take the latest com-  
645 plete timestamp of the included instances. For example: the aggregated task  
646 for case 9969 will have 13:45:17 as arrival time, 15:22:47 as start time and  
647 15:57:54 as completion time. For case 9974, the arrival, start and completion  
648 time correspond to 13:58:17, 15:57:54 and 16:28:57, respectively.

649 Finally, the third and last step uses the aggregated instances to determine  
650 whether *between-case batching* is present. Between-case batching determines  
651 whether a batching relationship also holds between the aggregated instances  
652 created in the second step (representing batches at the level of specific cases).  
653 To this end, similar sequential/concurrent time relationships are detected as  
654 the ones used for batch detection at the task-resource level. When between-  
655 case batching is also detected, a case-based batch subprocess is identified. In  
656 Table 8, a sequential relationship holds for the aggregated instances of cases  
657 9969 and 9974. Hence, a case-based sequential batch subprocess is found,  
658 which is highlighted by adding a shared batch subprocess identifier. The  
659 batch subprocess is also visualized in Figure 11.

660 The aforementioned three steps make it possible to determine whether  
661 case-based batching prevails for one particular task subsequence. However,  
662 a multitude of task subsequences can be present in the input file for phase

Table 8: Illustration of a sequential case-based batch subprocess (running example 1)

case id	task	resource	$T_{arrival}$	$T_{start}$	$T_{complete}$	task-resource batch id	task-resource batch type	batch subprocess id	batch subprocess type
...	...	...	...	...	...	...	...	...	...
9969	Study summary results	Lab technician June	14/01/2019 13:45:17	14/01/2019 15:22:47	14/01/2019 15:41:08	-	-	61	seq case-based
9969	Prepare report	Lab technician June	14/01/2019 15:41:08	14/01/2019 15:41:08	14/01/2019 15:57:54	-	-	61	seq case-based
9974	Study summary results	Lab technician June	14/01/2019 13:58:17	14/01/2019 15:57:54	14/01/2019 16:11:12	-	-	61	seq case-based
9974	Prepare report	Lab technician June	14/01/2019 16:11:12	14/01/2019 16:11:12	14/01/2019 16:28:57	-	-	61	seq case-based
...	...	...	...	...	...	...	...	...	...

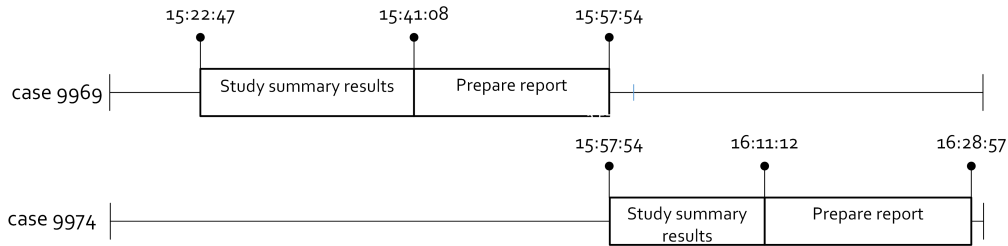


Figure 11: Visualization of a sequential case-based batch subprocess (running example 1)

663 3. Without loss of generality, two methods to identify candidate task subse-  
 664 quences are proposed, which are (1) enumeration and (2) sequence mining.  
 665 For each of these candidates, case-based batch detection will be performed  
 666 following the three steps outlined above.

667 When *enumeration* is used, all potential task subsequences prevailing in  
 668 the input subset for BAMA’s phase 3 are considered. This involves subse-  
 669 quences of all orders (i.e. number of tasks), going from two tasks to the length  
 670 of a trace included in the input subset. To avoid that very rare subsequences

671 also need to be checked for batching behavior, a threshold can be specified  
672 expressing the minimal occurrence frequency of a subsequence. Moreover,  
673 a filtering mechanism is integrated to avoid that higher-order subsequences  
674 are unnecessarily checked. Higher-order subsequences containing a particu-  
675 lar lower-order subsequence will not be checked when within-case batching  
676 is absent for this lower-order subsequence. The filtering mechanism builds  
677 upon the idea that within-case batching is a prerequisite for the detection  
678 of a case-based batch subprocess. Suppose, for example, that no within-case  
679 batching is detected for a task subsequence A - B. This implies that, for  
680 none of the cases in the input log, a batching relation exists for the case's  
681 instances related to tasks A and B. Hence, a case-based batch containing  
682 instances of A and B will not exist. Consequently, higher-order subsequences  
683 containing A - B such as A - B - C should not be checked as the prerequisite  
684 for case-based batching is not fulfilled for A - B.

685 Besides enumeration, existing *sequence mining* methods can also be used  
686 to identify candidate case-based subprocesses. In the implementation, the  
687 SPADE algorithm introduced by Zaki [39] has been incorporated. SPADE  
688 looks for frequent tasks sequences taking into account a user-defined min-  
689 imum support level. The implementation of BAMA hard-codes a SPADE  
690 parameter to ensure that only tasks that immediately follow each other are  
691 taken into consideration.

692 We now formally define the different batch subprocess types for case-based  
693 batching. First, requirements common for all case-based batch subprocess  
694 types are identified. Then, we define sequential, and concurrent case-based  
695 batching. Illustrations are provided to clarify the definitions, consisting of  
696 an example footprint in the batch-enriched task log and a visualization.

697 **Definition 8** (Base case-based batch conditions). *Let  $\psi \subseteq C$  represent a set*  
698 *of case identifiers. The base batch conditions common to all case-based batch*  
699 *types that a set of task instances  $T_\psi \subseteq T_{ex}$  needs to fulfill are:*

- 700
1. *all task instances refer to case identifiers in  $\psi$ , i.e.,  $\forall t \in T_\psi \#_c(t) \in \psi$ ;*
  2. *given  $1 \leq k \leq n-1$  and the sequence of task instances  $\langle t_1, \dots, t_k, \dots, t_n \rangle \in T_\psi^*$  with  $\#_{\tau_{start}}(t_k) \leq \#_{\tau_{start}}(t_{k+1})$  and  $t_k \neq t_{k+1}$  all task instances with the same case identifier are started before proceeding to task instances with another case identifier:*

$$\#_c(t_k) \neq \#_c(t_{k+1}) \implies \forall_{k+1 \leq j \leq n} (\#_c(t_j) \neq \#_c(t_k))$$

3. the same number and at least two instances for each case identifier are recorded:

$$\forall_{c_1, c_2 \in \psi} (|\{t \in T_\psi \mid \#_c(t) = c_1\}| = |\{t \in T_\psi \mid \#_c(t) = c_2\}| \geq 2);$$

4. the same set of task labels is involved in the batch:

$$\begin{aligned} & \forall_{c_1, c_2 \in \psi} (\{l \in L \mid \exists t \in T_\psi : \#_l(t) = l \wedge \#_c(t) = c_1\} \\ & = \{l \in L \mid \exists t \in T_\psi : \#_l(t) = l \wedge \#_c(t) = c_2\}); \end{aligned}$$

- 701 5. the same resource executed the full batch, i.e.,  $\forall_{t_a, t_b \in T_\psi} (\#_r(t_a) = \#_r(t_b))$ .

702 **Definition 9** (Sequential case-based batch). Let  $\psi \subseteq C$  represent a set  
703 of case identifiers. A sequential case-based batch is composed of a set of  
704 task instances  $T_\psi \subseteq T_{ex}$  fulfilling the base case-based batch conditions (cf.  
705 Definition 8) and additionally the following conditions:

- 706 1. for  $1 \leq k \leq n-1$  and the sequence of task instances  $\langle t_1, \dots, t_k, \dots, t_n \rangle \in$   
707  $T_\psi^*$  with  $\#_{\tau_{start}}(t_k) < \#_{\tau_{start}}(t_{k+1})$  the following conditions cumulatively  
708 hold:

- task instances referring to the same case identifier are at most separated by a time gap  $\gamma$ :

$$\begin{aligned} \#_c(t_k) = \#_c(t_{k+1}) & \implies (\#_{\tau_{start}}(t_{k+1}) - \#_{\tau_{complete}}(t_k)) \in [0, \gamma] \\ & \text{with } \gamma \geq 0, \end{aligned}$$

- the time gap between the completion of the last instance for a particular case identifier in  $\psi$  and the start of the first instance of the next case identifier in  $\psi$  cannot exceed the tolerated time gap  $\theta$ :

$$\begin{aligned} \#_c(t_k) \neq \#_c(t_{k+1}) & \implies (\#_{\tau_{start}}(t_{k+1}) - \#_{\tau_{complete}}(t_k)) \in [0, \theta] \\ & \text{with } \theta \geq 0. \end{aligned}$$

2. no other task instance performed by the same resource, which is not included in the sequential case-based batch, can be started or completed when processing the batched instances:

$$\begin{aligned} T_r = \{t \in (T_{ex} \setminus T_\psi) \mid & \exists_{t_a, t_b \in T_\psi} (\#_r(t_a) = \#_r(t) = \#_r(t_b)) \\ & \wedge (\#_{\tau_{start}}(t_a) \leq \#_{\tau_{start}}(t) \leq \#_{\tau_{complete}}(t_b)) \\ & \vee \#_{\tau_{start}}(t_a) \leq \#_{\tau_{complete}}(t) \leq \#_{\tau_{complete}}(t_b))\} \end{aligned}$$

709 we require  $T_r = \emptyset$ .

710 An illustration of a sequential case-based batch was already provided in  
 711 Table 8 and Figure 11 within the context of running example 1. Tasks ‘*Study*  
 712 *summary results*’ and ‘*Prepare report*’ are performed first for case 9969 in  
 713 a sequential way, and afterwards for case 9974. Moreover, the relationship  
 714 between the task instances of cases 9969 and 9974 also has a sequential  
 715 character. Consequently, a sequential case-based batch is detected.

716 **Definition 10** (Concurrent case-based batch). *Let  $\psi \subseteq C$  represent a set*  
 717 *of case identifiers. A concurrent case-based batch is composed of a set of*  
 718 *task instances  $T_\psi \subseteq T_{ex}$  fulfilling the base case-based batch conditions (cf.*  
 719 *Definition 8) and additionally the following conditions:*

1. *given  $1 \leq k \leq n-1$  and the sequence of task instances  $\langle t_1, \dots, t_k, \dots, t_n \rangle \in T_\psi^*$  with  $\#_{\tau_{start}}(t_k) \leq \#_{\tau_{start}}(t_{k+1})$  and  $t_k \neq t_{k+1}$  subsequent task instances are performed concurrently, but not perfectly in parallel:*

$$\begin{aligned} \#_{\tau_{start}}(t_k) \leq \#_{\tau_{start}}(t_{k+1}) < \#_{\tau_{complete}}(t_k) \\ \wedge (\#_{\tau_{start}}(t_k) \neq \#_{\tau_{start}}(t_{k+1}) \vee \\ \#_{\tau_{complete}}(t_k) \neq \#_{\tau_{complete}}(t_{k+1})) \end{aligned}$$

2. *no other task instance performed by the same resource, which is not included in the sequential case-based batch, can be started or completed when processing the batched instances;*

$$\begin{aligned} T_r = \{t \in (T_{ex} \setminus T_\psi) \mid \exists_{t_a, t_b \in T_\psi} (\#_r(t_a) = \#_r(t) = \#_r(t_b) \\ \wedge (\#_{\tau_{start}}(t_a) \leq \#_{\tau_{start}}(t) \leq \#_{\tau_{complete}}(t_b) \\ \vee \#_{\tau_{start}}(t_a) \leq \#_{\tau_{complete}}(t) \leq \#_{\tau_{complete}}(t_b))\} \end{aligned}$$

720 we require  $T_r = \emptyset$ .

721 To illustrate this last type of batch activity, consider Table 9 and Fig-  
 722 ure 12, which are derived from running example 2. For case 9123, tasks  
 723 ‘*Complete drug allergy form*’ and ‘*Perform blood test*’ are performed concu-  
 724 rrently. The same holds for case 9124. When aggregating both instances at  
 725 the case level, a concurrent relationship even exists between the cases due to  
 726 the time overlap between ‘*Perform blood test*’ for case 9123 and ‘*Complete*  
 727 *drug allergy form*’ for case 9124. Hence, a concurrent case-based subprocess  
 728 is detected.



Table 9: Illustration of a concurrent case-based batch subprocess (running example 2)

case id	task	resource	$\tau_{arrival}$	$\tau_{start}$	$\tau_{complete}$	task-resource batch id	task-resource batch type	batch subprocess id	batch subprocess type
...	...	...	...	...	...	...	...	...	...
9123	Complete drug allergy form	Nurse Kate	10/01/2019 15:02:41	10/01/2019 15:21:18	10/01/2019 15:33:09	-	-	21	conc case-based
9123	Perform blood test	Nurse Kate	10/01/2019 15:02:41	10/01/2019 15:31:30	10/01/2019 15:37:14	-	-	21	conc case-based
...	...	...	...	...	...	...	...	...	...
9124	Complete drug allergy form	Nurse Kate	10/01/2019 15:26:01	10/01/2019 15:35:24	10/01/2019 15:44:55	-	-	21	conc case-based
9124	Perform blood test	Nurse Kate	10/01/2019 15:26:01	10/01/2019 15:41:32	10/01/2019 15:48:09	-	-	21	conc case-based
...	...	...	...	...	...	...	...	...	...

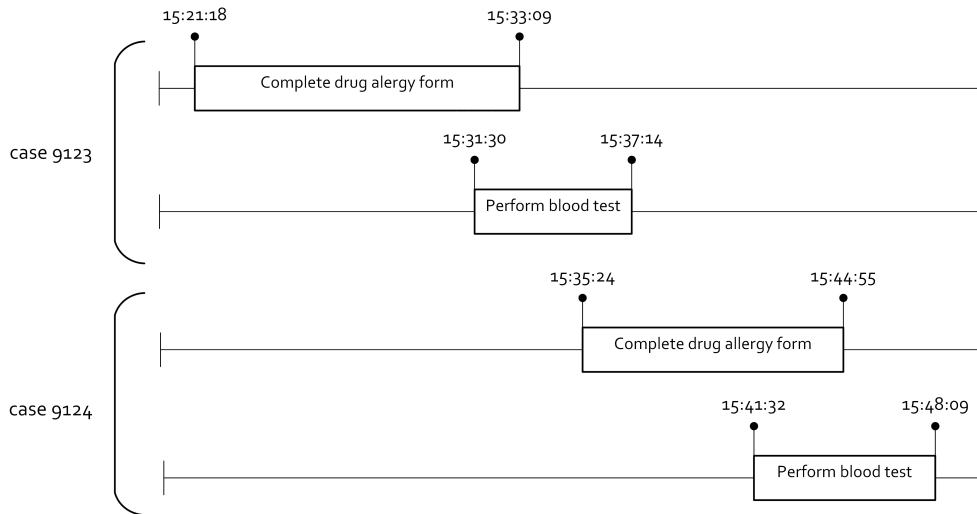


Figure 12: Visualization of a concurrent case-based batch subprocess (running example 2)

730 fill the requirements of a case-based sequential/concurrent batch subprocess.  
 731 This will be expressed in the output by assigning a shared batch subpro-  
 732 cess identifier to these instances. As a consequence, phase 3 will not add  
 733 new columns, but will complete the columns added in phase 2 (batch sub-  
 734 process identifier and batch subprocess type) with information on detected  
 735 case-based batch subprocesses.

#### 736 4.6. Batch-enriched task log creation

737 The output of BAMA is a batch-enriched task log, having a structure as  
 738 exemplified in Tables 5-9. It is obtained by merging the outputs of phases  
 739 2 and 3 of the algorithm. Compared to the task log, the batch-enriched  
 740 task log will contain four additional columns conveying batching information:  
 741 the task-resource batch identifier, the task-resource batch type, the batch  
 742 subprocess identifier, and the batch subprocess type. Hence, when applicable,  
 743 the algorithm will have grouped task instances in batches at the task-resource  
 744 level and identified batch subprocesses.

745 **Definition 11** (Batch-enriched task log). *Let  $T_{ex}$  represent an extended task*  
 746 *log and let  $T_{be}$  be a batch-enriched task log. To transform  $T_{ex}$  to  $T_{be}$ , two*  
 747 *mapping functions  $m_4$  and  $m_5$  are sequentially applied:*

- 748 •  $m_4 : T_{ex} \rightarrow T'_{ex}$  assigns a batch subprocess identifier  $b_{s,i}$  and a batch  
 749 subprocess type  $b_{s,t}$ , with  $b_{s,i} \in \mathbb{N}$  and  $b_{s,t} \in \{\text{par, seq task-based, conc}$   
 750  $\text{task-based}\}$ . The values of  $b_{s,i}$  and  $b_{s,t}$  are shared among the task in-  
 751 stances  $i$  which are part of the same parallel batch subprocess or sequen-  
 752 tial/concurrent task-based batch subprocess.
- 753 •  $m_5 : T'_{ex} \rightarrow T_{be}$  assigns a batch subprocess identifier  $b_{s,i}$  and a batch  
 754 subprocess type  $b_{s,t}$ , with  $b_{s,i} \in \mathbb{N}$  and  $b_{s,t} \in \{\text{seq case-based, conc case-}$   
 755  $\text{based}\}$ . The values of  $b_{s,i}$  and  $b_{s,t}$  are shared among the task instances  
 756  $i$  which are part of the same sequential/concurrent case-based batch  
 757 subprocess.

758 After the application of  $m_4$  and  $m_5$ , the batch-enriched task log  $T_{be}$  consists of  
 759 a set of batched task instances  $t'_b = (c, t, r, \tau_{arrival}, \tau_{start}, \tau_{complete}, b_{tr,i}, b_{tr,t}, b_{s,i},$   
 760  $b_{s,t})$ .

761 *4.7. Implementation*

762 A prototypical implementation of BAMA has been developed in R<sup>4</sup> and  
763 is publicly available at <https://github.com/nielsmartin/bama> or <https://doi.org/10.5281/zenodo.3653952>. R is a programming language provid-  
764 ing extensive functionalities for data manipulation and statistical analysis.  
765 The key packages that are used are `dplyr`<sup>5</sup> for data manipulations and sum-  
766 marisations, `lubridate`<sup>6</sup> to work with timestamps and `reshape` to convert  
767 the event log to a task log. When applying sequence mining during case-  
768 based batch detection, the implementation of the SPADE algorithm in the  
769 `arulesSequences`<sup>7</sup> package is used.

771 In order to use the implementation, an event log is required. This event  
772 log should take the form of a data frame, which is a rectangular data table for-  
773 mat in R in which columns contain variables and rows represent observations  
774 [40]. When an event log is available in the XES-format<sup>8</sup>, it can be converted  
775 to the correct input format by leveraging the R-package `xesreadR`<sup>9</sup>, which is  
776 part of the `bupaR` suite supporting process analysis in R [41].

777 Once the event log has been imported, a helper function is available to  
778 convert the event log to a task log. Afterwards, batch identification can be  
779 initiated. An integrated function enables the detection of batching behavior  
780 at both the task-resource and subprocess level. This function must be pa-  
781 rameterized by the user. The following parameters need to be specified to  
782 support the identification of batching behavior at the task-resource level:

- 783 • For sequential batch detection: a list containing the maximal tolerated  
784 time gaps between consecutive task instances. Each entry represents  
785 the time gap related to a particular task as, e.g., setup times might  
786 differ across tasks. A helper function is available to create this list.
- 787 • For sequential batch detection: a boolean indicating whether the event  
788 log contains the arrival time of a case at a task. When available, the  
789 arrival time is used to distinguish between sequential batching behavior  
790 and regular queue handling.

---

<sup>4</sup><https://www.r-project.org>

<sup>5</sup><https://CRAN.R-project.org/package=dplyr>

<sup>6</sup><https://CRAN.R-project.org/package=lubridate>

<sup>7</sup><https://CRAN.R-project.org/package=arulesSequences>

<sup>8</sup>XES is an XML-based standard for the exchange of event logs [2].

<sup>9</sup><https://CRAN.R-project.org/package=xesreadR>

- 791 • A boolean indicating whether timestamps are expressed numerically,  
792 instead of in a date-time format.
- 793 • The format of the timestamps in the event log (e.g. *yyyy-mm-dd*  
794 *hh:mm:ss*).

795 The following parameters need to be specified for the detection of sequen-  
796 tial/concurrent case-based batch subprocesses:

- 797 • An entry reflecting the way in which subsequences were generated, i.e.  
798 by means of enumeration or using sequence mining.
- 799 • A list containing the subsequences for which case-based sequential or  
800 concurrent batch subprocesses need to be detected. Helper functions  
801 are developed to create this list using enumeration or sequence mining  
802 using the SPADE algorithm.
- 803 • For sequential case-based batch detection: a maximal tolerated time  
804 gap between consecutive instances associated to a particular case (i.e.  
805 used to detect within-case batching).
- 806 • For sequential case-based batch detection: a maximal tolerated time  
807 gap between the (aggregated) within-case batches across consecutive  
808 cases (i.e. used to detect between-case batching).

#### 809 4.8. Metrics and visualizations

810 As follows from Section 4.6, the outcome of BAMA is a batch-enriched  
811 task log, which is a task log containing four additional columns with batch-  
812 ing information. As real-life task logs typically contain a large number of  
813 entries, a manual inspection of the batch-enriched task log is infeasible. Con-  
814 sequently, it is important that the batching information in the batch-enriched  
815 task log is presented in a concise and intuitive way in order to gain a rich un-  
816 derstanding in batching behavior. While a detailed overview on this matter  
817 is beyond the scope of this paper, the remainder of this subsection provides  
818 some pointers regarding (i) the potential of batching behavior metrics, and  
819 (ii) potential visualizations based on a batch-enriched task log.

820 A first way to convey batching insights in a concise way is the development  
821 of *batching behavior metrics*. The authors have discussed metrics for batch  
822 activities in previous works [6, 27]. A relevant metric is, for instance, the  
823 number of occurred batch activities as a percentage of the total number

824 of occurrences of the task(s). This allows insights into the actual usage  
825 of batching for certain tasks and subprocesses. Another metric can relate  
826 to the average or median batch size, providing insights in the number of  
827 cases which are typically included in a batch. Batching is usually applied  
828 to save time and costs, such that the average or median execution time  
829 and related costs of the batches can be compared to task executions which  
830 are not batched. When several resources can perform a particular batch  
831 activity, metrics can be calculated in general (i.e. disregarding the distinction  
832 between resources), or can be expressed at the level of individual resources.  
833 Positioning metrics at the resource level enables to analyze whether batching  
834 behavior is concentrated amongst a limited number of resources.

835 Besides numeric batching behavior metrics, the batch-enriched task log  
836 can also be used to *visualize batching behavior*. Graphs can take the metrics  
837 as a starting point and, e.g., show how the instances of a particular task  
838 are divided over the batch types. This can show the analyst how prevalent  
839 batching behavior is for a particular task. The evolution of particular metrics  
840 over time can also be visualized, which highlights whether batching behavior  
841 is related to particular times of day or days of the week. When considering  
842 a longer time horizon, it can also be observed whether batching behavior  
843 evolves over time. Next to graphs taking metrics as a starting point, the  
844 batching information in the batch-enriched task log can be projected on  
845 a Dotted Chart, in which the dots of batched task instances are colored,  
846 while non-batched instances are grey. Some examples of visualizations will  
847 be included in the evaluation of the proposed algorithm using real-life data  
848 (Section 5.2).

849 The batch-enriched task log can also be leveraged as an input for a visual  
850 analysis tool for batching behavior. Such a tool would enable an analyst to  
851 start from a high-level overview of batching behavior and interactively drill-  
852 down to study a specific batch or resource. Other features of such a tool  
853 could include the benchmark of resources, or the comparison between the  
854 characteristics of batched task instances and task instances which are not  
855 batched.

## 856 5. Evaluation

857 To evaluate BAMA, a twofold approach is followed: an evaluation using  
858 synthetic event logs and an evaluation using a real-life log. Both types of  
859 evaluation serve different goals. In Section 5.1, synthetic event logs are used

860 to evaluate the algorithm’s ability to rediscover known batches. Because it  
861 is known which task instances belonged to a batch in (simulated) reality,  
862 synthetic logs make it possible to demonstrate the effectiveness of BAMA  
863 to correctly detect batching behavior when it is present. Afterwards, in  
864 Section 5.2, the algorithm is applied to a real-life event log. To this end,  
865 historic data from a digital whiteboard system deployed in a surgical ward  
866 of a university hospital in Norway is used. Besides demonstrating BAMA’s  
867 applicability within a real-life setting, the insights that it can generate in  
868 such a real-life context are also highlighted.

## 869 5.1. Experiments with synthetic event logs

### 870 5.1.1. Experimental design

871 The goal of the evaluation using synthetic data is to determine whether  
872 the algorithm can rediscover known batches solely using an event log. To this  
873 end, synthetic event logs are created by simulating a BPMN process diagram  
874 using the extendable and open-source BPMN process simulator *Scylla* [42].  
875 Nine distinct models with varying degrees of complexity are included, as  
876 shown in Figure 13. We generated different variants of batch subprocesses  
877 (i) with a linear behavior (i.e., consisting of several tasks connected sequen-  
878 tially) including also a process model with two batch subprocesses (cf. #3  
879 in Figure 13), (ii) with parallel behavior through an AND gateway, (iii) with  
880 a choice through a XOR gateway, and (iv) with a combination of AND and  
881 XOR gateways.

882 For each process diagram, a separate event log is generated for each  
883 batch processing type, i.e. parallel, sequential/concurrent task-based, and  
884 sequential/concurrent case-based batching. As a consequence, 45 distinct  
885 event logs are generated. The number of simulated process instances per  
886 event log is set to 1,000. All synthetic event logs are publicly available  
887 at <https://github.com/nielsmartin/bama> or at <https://doi.org/10.5281/zenodo.3653952>.

889 For each synthetic event log, BAMA is applied using a maximum al-  
890 lowed time gap for sequential batching of zero seconds. Changing this value  
891 would have no impact on the results as the synthetic data originates from  
892 a simulated environment in which each task or batch activity has dedicated  
893 resources. Regarding the method to generate subsequences as candidates for  
894 case-based batching, *enumeration* is selected.

895 It should be stressed that, before feeding the event log to the algorithm,  
896 all information reflecting how batches are formed according to the simula-

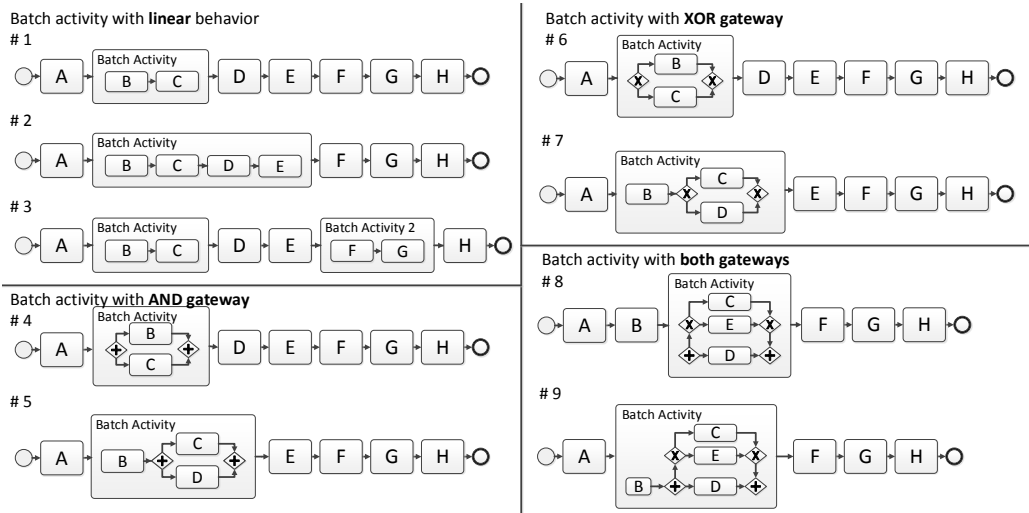


Figure 13: BPMN process diagrams used for the generation of synthetic event logs

897 tor’s batching logic is removed. Using this event log, completely agnostic  
 898 of the batches present in reality, the algorithm detects batches following the  
 899 procedure outlined in Section 4.

900 To evaluate the algorithm’s performance, the detected batches are compared  
 901 to the batches created during the simulation. The simulated environment  
 902 enables us to know which task instances are grouped together in a batch  
 903 according to the simulator’s batching logic. Hence, the algorithm should be  
 904 able to identify these batches from the synthetic log, providing a basis to eval-  
 905 uate whether the algorithm can identify different types of batching behavior  
 906 when we know that the input data contains such batching behavior. For per-  
 907 formance evaluation purposes, a task instance is considered to be incorrectly  
 908 assigned when it either (i) is included in a batch of the correct type, but in  
 909 a composition with instances different from the composition created by the  
 910 simulator’s batching logic, or (ii) is included in a wrong type of batch. Note  
 911 that the first criterion is rather rigorous as the composition of discovered  
 912 batches has to be completely correct. For instance: when BAMA rediscovers  
 913 a batch for all but one instance, the entire batch will be considered to be  
 914 incorrectly assigned.

### 915 5.1.2. Results

916 The use of synthetic data enables a direct comparison of the task instances  
 917 which are batched by the batching logic embedded in the simulator on the one

918 hand, and task instances which are batched according to BAMA’s output.  
919 To quantify the degree to which BAMA can correctly rediscover batches in  
920 synthetic data, the following three key output measures are calculated for  
921 each synthetic event log:

- 922 • **The percentage of correctly rediscovered batched instances.**  
923 This measure represents the percentage of task instances which are  
924 batched according to the simulator’s batching logic which BAMA cor-  
925 rectly rediscovers. In other words, this measure represents the percent-  
926 age of task instances batched according to the simulator’s batching  
927 logic, which are correctly marked as part of the correct batch by the  
928 algorithm.
  
- 929 • **The percentage of correctly rediscovered instances.** This mea-  
930 sure represents the percentage of task instances for which BAMA de-  
931 tects the correct batching behavior. Differently from the previous out-  
932 put measure, this measure also takes into account task instances which  
933 are not batched according to the simulator’s batching logic. As a con-  
934 sequence, this measure considers all task instances, while the previous  
935 measure focused on task instances which were batched according to the  
936 simulator’s batching logic.
  
- 937 • **The percentage of instances which are not batched according**  
938 **to the simulator’s batching logic, but which are reported as**  
939 **part of a sequential batch by BAMA.** This measure focuses on  
940 task instances which are not explicitly batched according to the sim-  
941 ulator’s batching logic. Despite the fact that these instances are not  
942 purposefully batched by the simulator, a sequential batching pattern  
943 might still be present because handling long queues can fulfill the re-  
944 quirements for sequential batching<sup>10</sup>. To quantify the extent to which  
945 this phenomenon is present, this measure highlights the percentage of  
946 instances which are not batched according to the simulator’s batching  
947 logic, but for which BAMA indicates that sequential batching occurs.

---

<sup>10</sup>Suppose that a large number of cases is waiting to be processed by a particular resource at a particular task. In that case, a group of queuing cases might be present at the task before the resource starts to process the first case in this group. As a consequence, under these conditions, sequential batching can be detected.



Table 10: Summary statistics on the output measures (excluding diagram #6)

<b>output measure</b>	<b>mean</b>	<b>sd</b>	<b>median</b>	<b>min</b>	<b>max</b>
% of correctly rediscovered batched instances	100.00	0.00	100.00	100.00	100.00
% of correctly rediscovered instances	90.12	2.65	90.97	83.33	91.67
% of instances not batched according to the simulator’s batching logic reported as part of sequential batch	16.19	6.84	13.49	11.21	33.33

948 Batch identification on all 45 synthetic event log is executed on a standard  
 949 laptop computer. The average runtime was well below 10 seconds. The  
 950 detailed results for each synthetic log are included in Table A.11 in Appendix  
 951 A. From these results, it follows that the algorithm will correctly rediscover all  
 952 batching behavior which is introduced according to the simulator’s batching  
 953 logic. The only exception is process diagram #6, where the batch subprocess  
 954 is not rediscovered as the model consists of two tasks in an XOR-construct.  
 955 As either activity B or C will be observed in a trace, the algorithm will not  
 956 discover that both tasks are part of a batch subprocess. This observation  
 957 will be revisited in the discussion (Section 6).

958 When disregarding diagram #6, which is a clear outlier in terms of perfor-  
 959 mance, Table 10 provides summary statistics on the three output measures.  
 960 The table shows that all task instances which the simulator batches accord-  
 961 ing to its batching logic are correctly rediscovered by the algorithm (100%  
 962 rediscovery). When instances which are not batched according to the simu-  
 963 lator’s batching logic are also taken into account, the algorithm classifies, on  
 964 average, 90.12% of the instances correctly. Given the fact that all instances  
 965 that the simulator batches according to its batching logic are correctly re-  
 966 discovered, this last result indicates that BAMA reports some task instances  
 967 which are not batched by the simulator as being part of a batch. This is  
 968 supported by the last output measure in Table 10: on average 16.19% of the  
 969 instances which are not batched according to the simulator’s batching logic  
 970 are marked as being part of a sequential batch by BAMA. This latter obser-  
 971 vation explains the entire deviation between the designed batching behavior  
 972 in the simulator and our algorithm’s output.

973 From what said, the following two observations follow: (i) all batched  
 974 task instances are correctly rediscovered for all process diagrams besides di-  
 975 agram #6, and (ii) the identified batch behavior is not correctly discovered  
 976 for all task instances because BAMA detects sequential batching behavior  
 977 which is not introduced according to the simulator’s batching logic. These

978 observations will be discussed in Section 6.

979 As highlighted in the experimental design, *enumeration* was used as the  
980 method to generate subsequences which are candidate for case-based batch-  
981 ing. For the sake of completeness, it should be noted that rerunning all  
982 experiments with *sequence mining* as a subsequence identification method  
983 generated identical results. Differences in terms of runtime were minimal  
984 and remained, on average, well below 10 seconds for each synthetic log using  
985 a standard laptop computer.

## 986 5.2. Discovering Batch Activities in Usage of a Hospital Ward

987 BAMA is also applied to a real-life event log obtained from a digital  
988 whiteboard system that is deployed in a surgical ward of an university hospi-  
989 tal in Norway. The event log has been subject to a previous study in which  
990 higher level activities were recognized from the low-level events recording  
991 every change in the system [43]. The event log is a good candidate for eval-  
992 uating BAMA as it is expected to contain some batching behavior in how  
993 information is entered for patients. Indeed, batching behavior was already  
994 presumed in Mannhardt et al. [43] albeit solely based on a preliminary visual  
995 analysis using a Dotted Chart [44]. Specifically, the hypothesis regarding  
996 batching was that nurses would often only use the whiteboard to update  
997 certain information around the times of a shift change.

998 The whiteboard system was introduced as a light-weight support system  
999 for coordination and collaboration among nurses and is meant to support the  
1000 daily work of nurses in the hospital [46]. Figure 14 shows the main screen of  
1001 the whiteboard in which each row corresponds to all the information entered  
1002 for a single patient. Each patient is assigned to a responsible nurse and that  
1003 nurse can update information for that patient in the whiteboard, e.g., the  
1004 planned treatment, necessary reports, or the planned discharge date. Next  
1005 to medical information, the whiteboard is also linked to a call signal system,  
1006 which records when patients raise an alarm or nurses attend the patient in  
1007 their room. Finally, the whiteboard system also records when the responsible  
1008 nurse for a patient changes. Thus, resource information is available for all  
1009 events.

### 1010 5.2.1. Experimental design

1011 We obtained an event log with 8,487 cases tracking the updates made  
1012 on the whiteboard system for individual patients. In total, there are 298,636  
1013 events recorded. The events are recorded at fine granularity, i.e. every change

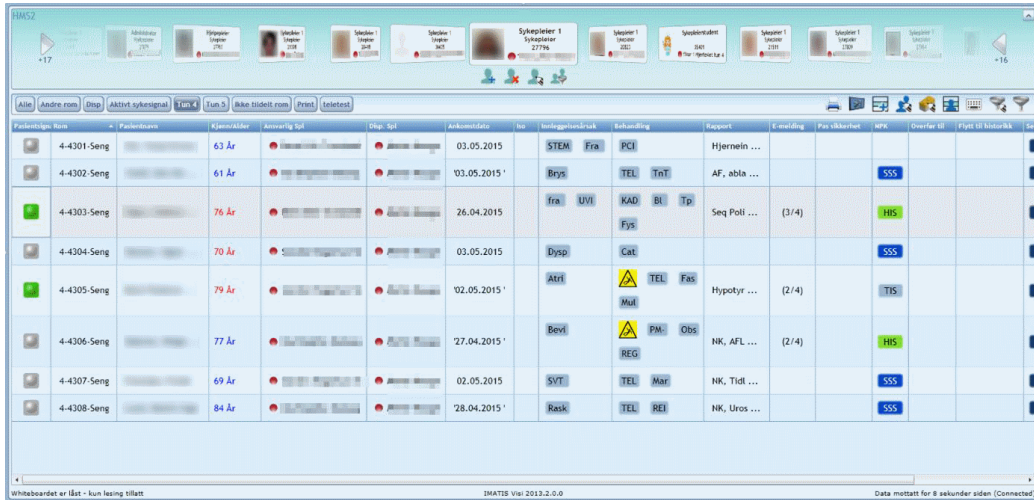


Figure 14: Whiteboard system used in the hospital as described in [45]. Each row contains information about a patient can be updated directly through a touch screen interface. The event log of the system contains events for all the updates made grouped by patient. Note that some entries are purposefully blurred for anonymization purposes.

1014 of a cell in the whiteboard yields an event, and carries only a single times-  
 1015 tamp as payload. We used the abstraction method and patterns described  
 1016 in Mannhardt et al. [43] to derive a high-level event log with 6,455 cases<sup>11</sup>  
 1017 and 70,936 instances of 14 distinct activities that have a defined *start* and  
 1018 *completion* timestamp. The activities refer to:

- 1019 • pre-announcement, registration, and transfer of the patient;
- 1020 • five different usage patterns of the nurse call signal system;
- 1021 • updates of treatment and diagnosis information as well as generic re-  
 1022 ports on the patient;
- 1023 • handover between nurses.

1024 Further details on the semantics of all activities are provided in Mannhardt  
 1025 et al. [43].

<sup>11</sup>For some cases the abstraction method did result in empty traces as no activity was recognized. These cases were filtered out.

1026 BAMA is applied to the prepared event log. We used SPADE to identify  
 1027 frequent subsequences as exhaustively computing all subsequences took  
 1028 about 10 minutes and returned more than 130,000 subsequences, which would  
 1029 result in a very long computation time when applying BAMA. By varying  
 1030 the minimum support level parameter between 0.005 and 0.4, we investigated  
 1031 the trade-off on the frequency of detected case-based batching and found 0.01  
 1032 to be giving the best result within a computation time of 3 minutes. The  
 1033 tolerated time gap parameters of BAMA were set to 3 minutes based on the  
 1034 domain knowledge that some activities may incur additional work beyond  
 1035 what is captured in the timeframe between the *start* and *complete* time of  
 1036 an activity instance. Again, we investigated the impact of the parameter by  
 1037 varying it from 30 seconds to 12 minutes. As expected, more batch behavior  
 1038 is detected when increasing the parameter value. BAMA detects 15.7% more  
 1039 batching when using 12 minutes compared to 30 seconds. Finally, in absence  
 1040 of more precise information in the event log, we assumed the arrival time of  
 1041 tasks to be 5 minutes before their start.

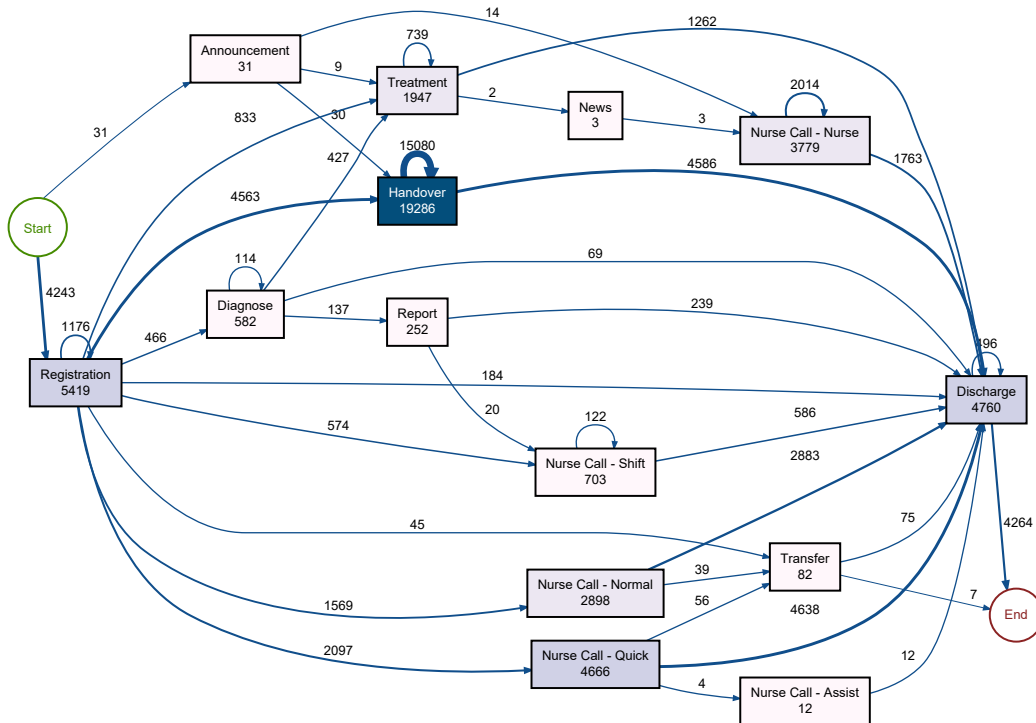


Figure 15: Process map discovered with Flexible Heuristics Miner (FHM) on the white-board event log showing the main activities supported by the whiteboard system.

1042 Figure 15 gives an overview of the process behavior for all completed  
1043 cases, i.e., those starting with either *‘Registration’* or *‘Announcement’* and  
1044 ending with *‘Discharge’* or *‘Transfer’*. The process map shows causal depen-  
1045 dencies discovered with the Flexible Heuristics Miner (FHM) together with  
1046 projected frequencies from the event log. As expected there is little struc-  
1047 ture since most of the activities can be performed in any order. Some of  
1048 the activities can be repeated such as *‘Handover’*, *‘Nurse Call - Nurse’* and  
1049 *‘Nurse Call - Shift’*, as well as activities which update patient information  
1050 such as *‘Diagnose’* and *‘Treatment’*. The former three activities are initiated  
1051 by the nurse, who visits the patient’s room without a request from the pa-  
1052 tient. Finally, both *‘Registration’* and *‘Discharge’* are repeated, which can  
1053 be attributed to technical issues as the source system duplicates these events.

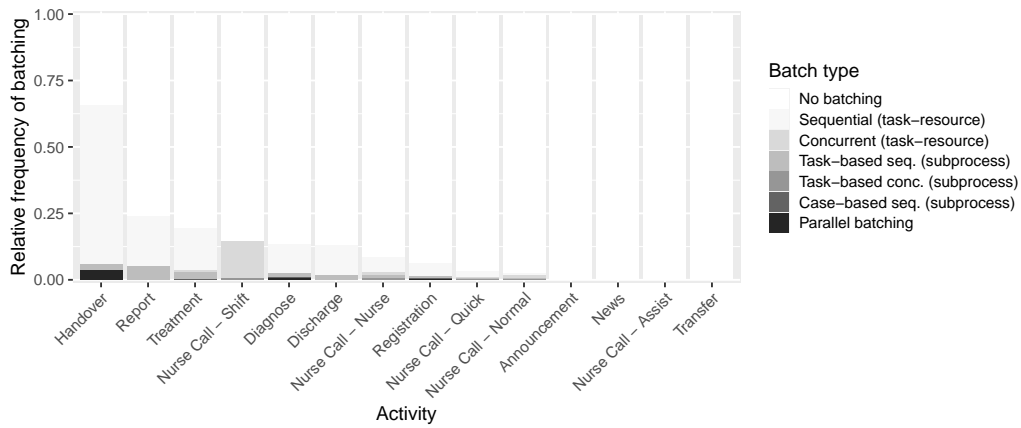
### 1054 5.2.2. Results

1055 The BAMA algorithm detected 30,530 task instances that were batched  
1056 in the digital whiteboard system. The execution time was about 3 minutes  
1057 on a standard laptop computer. Figure 16 shows an overview of the relative  
1058 frequencies with which batching behavior is detected for the different activ-  
1059 ities. This overview includes both batched instances at the task-resource  
1060 level, which could already be detected by the state-of-the-art technique de-  
1061 scribed in Martin et al. [6], as well as batch subprocesses, which are the core  
1062 contribution of this work.

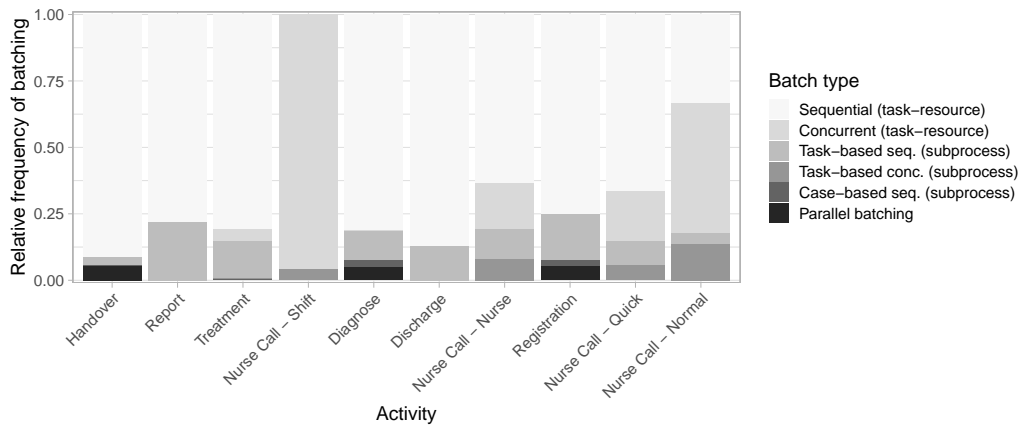
1063 The most frequent task related to batching behavior is *‘Handover’*. This  
1064 is not surprising as handover of work normally takes place at the end of a  
1065 shift and the change in responsibility for each patient is registered on the  
1066 whiteboard. So, it is to be expected that this task is executed by nurses as a  
1067 sequential batch. Beyond this obvious batching behavior, BAMA identified  
1068 batching behavior in more than 20% of the task instances for the activities  
1069 *‘Report’*, *‘Treatment’*, *‘Diagnose’*, and *‘Discharge’* as well as two types of  
1070 interactions with the call signal system: *‘Nurse Call - Shift’* and *‘Nurse Call*  
1071 *- Nurse’*. In the latter two tasks, nurses use the call signal to indicate their  
1072 position in the ward.

1073 These results partially confirm the batching hypothesis that was formu-  
1074 lated in Myrstad [45]. However, in contrast to the basic analysis with a Dot-  
1075 ted Chart visualization in Myrstad [45], the application of BAMA accurately  
1076 quantifies the amount of batching behavior, which is less than expected.

1077 Next, we looked at the batching behavior at the subprocess level that  
1078 was identified by BAMA. Figure 17 shows how frequently a task appears to



(a)



(b)

Figure 16: Distribution of the batch types (cf. Figure 4) identified for task instances. In (a), instances for which no batching behavior is detected are included (*No batching*). These instances are excluded in (b).

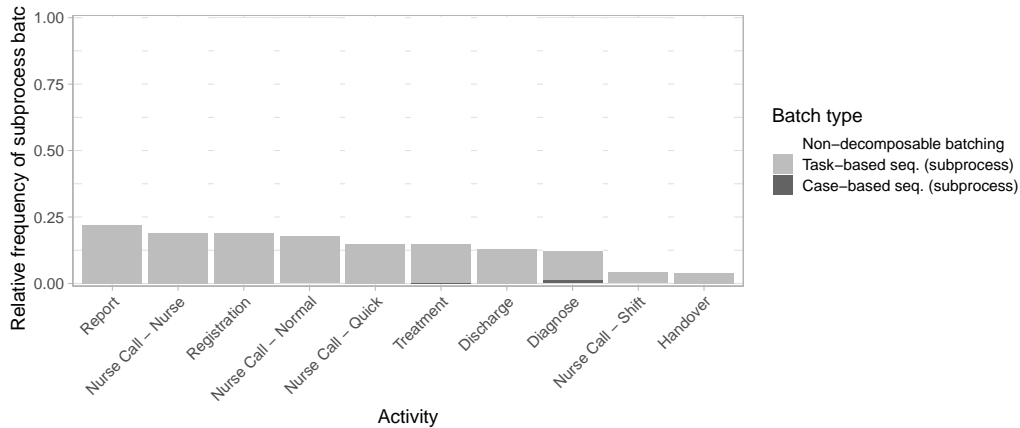


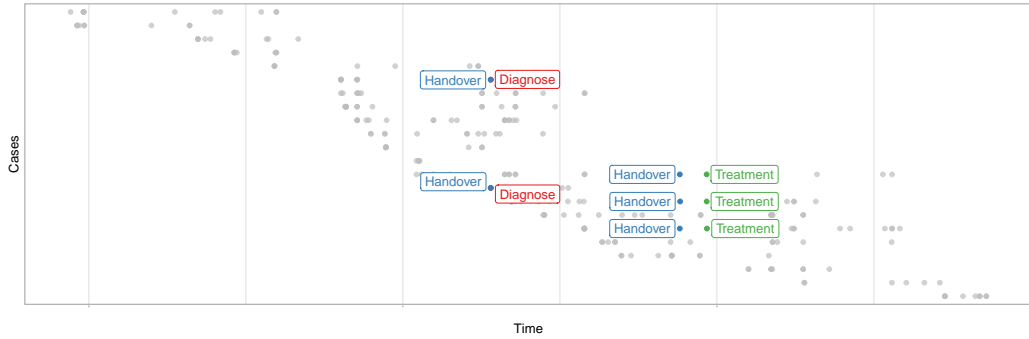
Figure 17: Distribution of type of batching at the subprocess level compared based on the overall batching.

1079 be batched in a subprocess in relation to the overall batching behavior that  
 1080 was identified for that task, i.e., the number of activity instances in sub-  
 1081 process batching compared to the overall number of instances performed in  
 1082 batches. Overall, subprocess batching was observed infrequently in the event  
 1083 log. For example, in about 12% of the cases, the tasks ‘*Nurse Call - Normal*’  
 1084 and ‘*Report*’ were identified as being part of sequential subprocess batching.  
 1085 Concurrent subprocess batching is almost entirely absent, which is to be ex-  
 1086 pected since there is normally only one nurse responsible for the patients on a  
 1087 ward. We did not discover parallel batch subprocesses. In total, BAMA  
 1088 discovered 1,220 task instances that were executed in subprocess batching  
 1089 out of which 1,194 instances were part of task-based batch subprocesses and  
 1090 26 were part of case-based batch subprocesses.

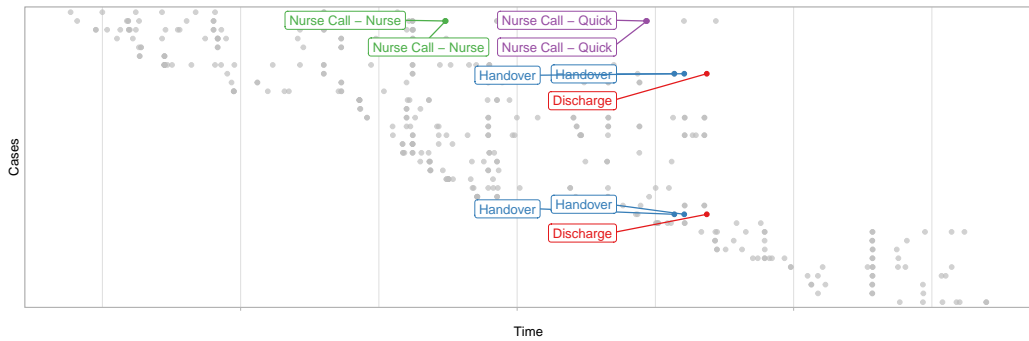
1091 In Figure 18, we grouped the detected batch subprocesses by the tasks  
 1092 involved and indicate their ordering with the symbol ‘ $\rightarrow$ ’ to investigate the  
 1093 batching behavior in more detail. Most of the subprocesses include the task  
 1094 ‘*Handover*’ batched together with tasks ‘*Discharge*’, ‘*Treatment*’, ‘*Nurse Call*  
 1095 - *Nurse*’, ‘*Report*’, and ‘*Registration*’. Whereas batching of the ‘*Handover*’  
 1096 task is to be expected, entering information about the treatment for patients  
 1097 (‘*Treatment*’) or entering information in the report field of the whiteboard  
 1098 (‘*Report*’) does not necessarily need to be batched. The whiteboard is sup-  
 1099 posed to be used continuously during the shift to always have the latest  
 1100 information available. However, as already investigated in Myrstad [45], it is  
 1101 often only entered afterwards.



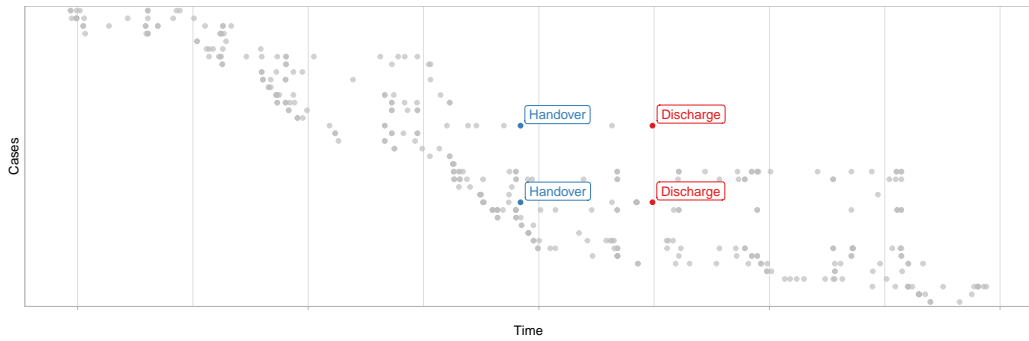




(a)



(b)



(c)

Figure 19: Dotted chart of events recorded in the observation ward with those events highlighted that have been detected by BAMA as task-based subprocesses within the course of a single day. All batch subprocesses are sequential, i.e., the task instances did not occur in parallel, which is not visible due to the coarse time scale. The exact time is not revealed for privacy reasons.

1123 tasks in an XOR-construct. Consequently, either task B or task C will be  
1124 executed for a batch. Hence, the algorithm will identify batching behavior at  
1125 the task-resource level for B or C instead of a batch subprocess consisting of  
1126 B and C. To circumvent this limitation, BAMA’s output can be projected on  
1127 a process model generated using an existing control-flow discovery algorithm.  
1128 This will highlight the relationship between these batch tasks and, hence, the  
1129 fact that a batch subprocess is formed.

1130 Another observation from the evaluation using synthetic data is that  
1131 BAMA detects more batching behavior than the batches which were pur-  
1132 posefully introduced by the simulator. More specifically, some task instances  
1133 which are not batched according to the simulator’s batching logic are marked  
1134 by BAMA as part of a sequential batch (at the task-resource level). Even  
1135 though this behavior is not introduced by the simulator’s batching logic, it  
1136 presents valid batching behavior as long queues which are sequentially pro-  
1137 cessed can fulfill the criteria for sequential batching. In particular, sequential  
1138 batching requires that a group of cases is queuing for the resource at a task  
1139 before this resource starts processing the first case within this group. When  
1140 queues are long, this condition can be fulfilled. Against this background, the  
1141 detection of sequential batching behavior by BAMA does not constitute a  
1142 performance issue of the algorithm.

1143 The application of the algorithm on real-life data shows that BAMA iden-  
1144 tifies batching both at the task-resource level and at the subprocess level.  
1145 Overall, the application of BAMA generated more profound insights into  
1146 batching behavior than previous analyses of the same data with standard  
1147 methods, such as a Dotted Chart analysis. In particular, BAMA can be  
1148 used to exactly quantify the fraction of batch executions for certain parts  
1149 of the process. Compared to the previous analysis in Mannhardt et al. [43]  
1150 and Myrstad [45], batching of certain subprocesses (e.g., ‘*Report*’ followed  
1151 by ‘*Handover*’) was confirmed, but at a lower frequency than expected. The  
1152 generated insights provides support to the hospital to investigate the op-  
1153 erational use of the whiteboard system. For instance: the observation that  
1154 particular reporting tasks are batched together with handover tasks indicates  
1155 that the whiteboard is not always used continuously during the shift.

1156 While the synthetic logs and the real-life log in Section 5 originate from a  
1157 single process, BAMA can also identify batching behavior in a multi-process  
1158 context. This is due to the fact that the algorithm focuses on how resources  
1159 perform tasks. Consequently, when tasks of several processes are included in  
1160 the event log, BAMA can also detect batch subprocesses containing tasks of

1161 different processes.

## 1162 *6.2. Limitations*

1163 Besides the contributions of this paper, its limitations also need to be  
1164 recognized. Firstly, the algorithm will identify a batch subprocess solely  
1165 consisting of a set of tasks in an XOR-construct as batching behavior at the  
1166 task-resource level. As outlined above, this limitation can be circumvented  
1167 by projecting the algorithm’s output on the output of a control-flow discovery  
1168 algorithm.

1169 Secondly, the algorithm imposes some requirements on the event log used  
1170 as an input. Foremost, the start and completion of a task instance need to be  
1171 recorded, which might not be the case in particular real-life event logs. When  
1172 events are recorded at a fine granularity, this issue can be tackled by applying  
1173 abstraction methods, as discussed in the real-life data evaluation. Besides  
1174 start and completion times, the arrival time of a case at an activity is used  
1175 to distinguish sequential batching from regular queue handling. However, it  
1176 should be noted that the absence of an arrival time or a suitable proxy does  
1177 not impede the application of BAMA. When the arrival time is not available,  
1178 the conditions for sequential batching are less stringent as the requirement  
1179 that all cases need to be present before the processing of a batch starts cannot  
1180 be enforced.

1181 Finally, the algorithm centers around the identification of batching behav-  
1182 ior, but does not focus on the operational effects of such behavior. However,  
1183 the algorithm identifies a wide variety of batching behavior, both at the  
1184 task-resource level and at the subprocess level. This provides a solid basis  
1185 to investigate why particular cases are batched and whether this is desirable  
1186 from a performance perspective.

## 1187 **7. Conclusion**

1188 This paper presents the Batch Activity Mining Algorithm (BAMA), which  
1189 is a novel algorithm to automatically detect batching behavior from an event  
1190 log. It extends prior research on this matter by enabling the discovery of  
1191 both batch tasks and batch subprocesses. The evaluation, both on synthetic  
1192 and real-life data, demonstrates the algorithm’s ability to identify batches  
1193 in an event log. BAMA’s output provides organizations with quantitative  
1194 insights in the occurrence of a wide variety of batching behaviors. Compared  
1195 to a mere visual analysis, e.g. using Dotted Charts, the presented algorithm

1196 provides a more structured approach to identify batching behavior. Batching  
1197 can have a positive impact (e.g. a reduction in the number of setups), as well  
1198 as a negative impact (e.g. an increase in waiting times for some cases) on  
1199 the performance of a process. The algorithm offers detailed insights into the  
1200 recorded batch behavior and can be helpful to explain the observed process  
1201 duration in more detail. As the algorithm focuses on the relation between  
1202 tasks and resources, it could also detect batch behavior over multiple pro-  
1203 cesses if tasks of the different processes are included in one event log.

1204 Building upon the work presented in this paper, several interesting di-  
1205 rections for future research can be distinguished. Firstly, BAMA can be  
1206 extended with a set of aggregated metrics and visualization functions. Cur-  
1207 rently, the algorithm focuses on the identification of batch tasks and sub-  
1208 processes. Taking this output as a starting point, future developments could  
1209 create a framework to, e.g., interactively analyze batching behavior. Some  
1210 pointers were already included in Section 4.8. Secondly, future work can link  
1211 the identified batches to its operational effects (e.g. its impact on waiting  
1212 times from a customer’s perspective) to determine whether a particular type  
1213 of batching behavior is desirable from the organizational perspective. More-  
1214 over, the organization could determine whether particular persons or teams  
1215 are more inclined to exhibit batching behavior. Finally, it is worthwhile to  
1216 investigate whether inductive batching insights can be embedded in a predic-  
1217 tive process monitoring framework. This would imply that knowledge about  
1218 batching behavior from the past would be used as one of the predictors to  
1219 estimate, e.g., the remaining time required to finish a case.

1220 **Acknowledgement.** We would like to thank Leon Bein (Master student at  
1221 Hasso Plattner Institute) for extending the simulator Scylla and for support-  
1222 ing the generation of the syntactic event logs. We would also like to sincerely  
1223 thank the reviewers for their constructive feedback during the review process.

## 1224 **Appendix A. Detailed output measures of the evaluation on syn-** 1225 **thetic data**

- 1226 1. Process diagram #1: simple var1 - one batch activity with two tasks
- 1227 2. Process diagram #2: simple var2 - one batch activity with four tasks
- 1228 3. Process diagram #3: simple var3 - two batch activities, each consisting  
1229 of two tasks
- 1230 4. Process diagram #4: AND var 1 - one batch activity with two tasks

- 1231 5. Process diagram #5: AND var 2 - one batch activity with three tasks
- 1232 6. Process diagram #6: XOR var 1 - one batch activity with two tasks in
- 1233 XOR construct
- 1234 7. Process diagram #7: XOR var 2 - one batch activity with one 'fixed'
- 1235 task and two tasks in XOR construct
- 1236 8. Process diagram #8: MIX var 1 - one batch activity with AND followed
- 1237 by XOR in one branch
- 1238 9. Process diagram #9: MIX var 2 - one batch activity with one 'fixed'
- 1239 task, followed by an AND with an XOR in one branch

Table A.11: Detailed output measures of evaluation on synthetic data

log number	process diagram	batch processing type*	% of correctly rediscovered instances	% of correctly rediscovered instances	% of instances not batched according to the simulator's batching logic reported as seq. batch
1	1	par.	100.00	91.52	11.31
2	1	s.t.b.	100.00	91.49	11.34
3	1	s.c.b.	100.00	91.59	11.21
4	1	c.t.b.	100.00	91.49	11.34
5	1	c.c.b.	100.00	91.59	11.21
6	2	par.	100.00	91.67	16.67
7	2	s.t.b.	100.00	91.67	16.67
8	2	s.c.b.	100.00	91.67	16.67
9	2	c.t.b.	100.00	91.67	16.67
10	2	c.c.b.	100.00	91.67	16.67
11	3	par.	100.00	83.33	33.33
12	3	s.t.b.	100.00	83.33	33.33
13	3	s.c.b.	100.00	83.33	33.33
14	3	c.t.b.	100.00	83.33	33.33
15	3	c.c.b.	100.00	83.33	33.33
16	4	par.	100.00	91.47	11.38
17	4	s.t.b.	100.00	91.59	11.21
18	4	s.c.b.	100.00	91.54	11.28
19	4	c.t.b.	100.00	91.59	11.21
20	4	c.c.b.	100.00	91.54	11.28
21	5	par.	100.00	91.59	13.45
22	5	s.t.b.	100.00	91.57	13.49
23	5	s.c.b.	100.00	91.57	13.49
24	5	c.t.b.	100.00	91.57	13.49
25	5	c.c.b.	100.00	91.57	13.49
26	6	par.	0.00	76.05	11.28
27	6	s.t.b.	0.00	76.02	11.31
28	6	s.c.b.	0.00	76.16	11.14
29	6	c.t.b.	0.00	76.02	11.31
30	6	c.c.b.	0.00	76.16	11.14
31	7	par.	100.00	90.25	13.65
32	7	s.t.b.	100.00	90.36	13.49
33	7	s.c.b.	100.00	90.39	13.45
34	7	c.t.b.	100.00	90.36	13.49
35	7	c.c.b.	100.00	90.39	13.45
36	8	par.	100.00	90.48	13.33
37	8	s.t.b.	100.00	90.48	13.33
38	8	s.c.b.	100.00	90.48	13.33
39	8	c.t.b.	100.00	90.48	13.33
40	8	c.c.b.	100.00	90.48	13.33
41	9	par.	100.00	90.48	16.67
42	9	s.t.b.	100.00	90.48	16.67
43	9	s.c.b.	100.00	90.48	16.67
44	9	c.t.b.	100.00	90.48	16.67
45	9	c.c.b.	100.00	90.48	16.67

\* *par.*: parallel, *s.t.b.*: sequential task-based, *s.c.b.*: sequential case-based, *c.t.b.*: concurrent task-based, *c.c.b.*: concurrent case-based

- 1240 [1] M. Weske, Business process management: concepts, languages, archi-  
1241 tectures, Springer-Verlag Berlin Heidelberg, 2012.
- 1242 [2] W. M. P. van der Aalst, Process mining: data science in action, Springer,  
1243 Heidelberg, 2016.
- 1244 [3] L. Pufahl, M. Weske, Requirements framework for batch processing in  
1245 business processes, Lecture Notes in Business Information Processing  
1246 287 (2017) 85–100.
- 1247 [4] C. N. Potts, M. Y. Kovalyov, Scheduling with batching: a review,  
1248 European journal of operational research 120 (2000) 228–249.
- 1249 [5] J. Medhi, Stochastic models in queueing theory, Academic Press, 2002.
- 1250 [6] N. Martin, M. Swennen, B. Depaire, M. Jans, A. Caris, K. Vanhoof,  
1251 Retrieving batch organisation of work insights from event logs, Decision  
1252 Support Systems 100 (2017) 119–128.
- 1253 [7] Y. Wen, Z. Chen, J. Liu, J. Chen, Mining batch processing workflow  
1254 models from event logs, Concurrency and Computation: Practice and  
1255 Experience 25 (2013) 1928–1942.
- 1256 [8] J. Nakatumba, Resource-aware business process management: analysis  
1257 and support, Ph.D. thesis, Eindhoven University of Technology, 2013.
- 1258 [9] Y. Liu, L. Zhang, J. Wang, Mining workflow event log to facilitate par-  
1259 allel work item sharing among human resources, International Journal  
1260 of Computer Integrated Manufacturing 24 (2011) 864–877.
- 1261 [10] L. Pufahl, A. Meyer, M. Weske, Batch regions: process instance syn-  
1262 chronization based on data, Proceedings of the 2014 IEEE International  
1263 Enterprise Distributed Object Computing Conference (2014) 150–159.
- 1264 [11] N. Slack, S. Chambers, R. Johnston, Operations and process manage-  
1265 ment: principles and practice for strategic impact, Pearson Education,  
1266 2009.
- 1267 [12] S. Henn, S. Koch, G. Wäscher, Order batching in order picking ware-  
1268 houses: a survey of solution approaches, Springer, 2012.

- 1269 [13] S. Henn, G. Wäscher, Tabu search heuristics for the order batching prob-  
1270 lem in manual order picking systems, *European Journal of Operational*  
1271 *Research* 222 (2012) 484–494.
- 1272 [14] S. Hong, A. L. Johnson, B. A. Peters, Batch picking in narrow-aisle  
1273 order picking systems with consideration for picker blocking, *European*  
1274 *Journal of Operational Research* 221 (2012) 557–570.
- 1275 [15] S. Hong, Y. Kim, A route-selecting order batching model with the S-  
1276 shape routes in a parallel-aisle order picking system, *European Journal*  
1277 *of Operational Research* 257 (2017) 185–196.
- 1278 [16] J. K. Higginson, J. H. Bookbinder, Policy recommendations for a  
1279 shipment-consolidation program, *Journal of Business Logistics* 15  
1280 (1994).
- 1281 [17] S. Cetinkaya, J. H. Bookbinder, Stochastic models for the dispatch of  
1282 consolidated shipments, *Transportation Research Part B: Methodolog-*  
1283 *ical* 37 (2003) 747–768.
- 1284 [18] F. Mutlu, S. i. l. Cetinkaya, J. H. Bookbinder, An analytical model for  
1285 computing the optimal time-and-quantity-based policy for consolidated  
1286 shipments, *IIE Transactions* 42 (2010) 367–377.
- 1287 [19] L. Tang, G. Wang, Decision support system for the batching problems  
1288 of steelmaking and continuous-casting production, *Omega* 36 (2008)  
1289 976–991.
- 1290 [20] C. N. Potts, M. Y. Kovalyov, Scheduling with batching: a review,  
1291 *European Journal of Operational Research* 120 (2000) 228–249.
- 1292 [21] M. Mathirajan, A. I. Sivakumar, A literature review, classification and  
1293 simple meta-analysis on scheduling of batch processors in semiconductor,  
1294 *The International Journal of Advanced Manufacturing Technology* 29  
1295 (2006) 990–1001.
- 1296 [22] M. F. Neuts, A general class of bulk queues with Poisson input, *The*  
1297 *Annals of Mathematical Statistics* 38 (1967) 759–770.
- 1298 [23] K. Sikdar, U. C. Gupta, Analytic and numerical aspects of batch ser-  
1299 vice queues with single vacation, *Computers & Operations Research* 32  
1300 (2005) 943–966.



- 1301 [24] C. Natschläger, A. Bögl, V. Geist, M. Biró, Optimizing resource uti-  
1302 lization by combining activities across process instances, in: EuroSPI,  
1303 Springer, pp. 155–167.
- 1304 [25] J. Pflug, S. Rinderle-Ma, Application of dynamic instance queuing to ac-  
1305 tivity sequences in cooperative business process scenarios, *International*  
1306 *Journal of Cooperative Information Systems* (2016) 1650002.
- 1307 [26] L. Pufahl, M. Weske, Batch activity: enhancing business process mod-  
1308 eling and enactment with batch processing, *Computing* (2019) 1–25.
- 1309 [27] L. Pufahl, E. Bazhenova, M. Weske, Evaluating the performance of a  
1310 batch activity in process models, *Lecture Notes in Business Information*  
1311 *Processing* 202 (2014) 277–290.
- 1312 [28] N. Martin, A. Solti, J. Mendling, B. Depaire, A. Caris, Mining batch  
1313 activation rules from event logs, *IEEE Transactions on Services Com-*  
1314 *puting* (2019).
- 1315 [29] A. Senderovich, M. Weidlich, A. Gal, A. Mandelbaum, Queue min-  
1316 ing: predicting delays in service processes, *Lecture Notes in Computer*  
1317 *Science* 8484 (2014) 42–57.
- 1318 [30] I. Weber, M. Farshchi, J. Mendling, J.-G. Schneider, Mining processes  
1319 with multi-instantiation, in: *Proceedings of the 30th Annual ACM*  
1320 *Symposium on Applied Computing*, pp. 1231–1237.
- 1321 [31] OMG, Notation BPMN version 2.0, OMG Specification, Object Man-  
1322 agement Group (2011).
- 1323 [32] A. Senderovich, C. Di Francescomarino, C. Ghidini, K. Jorbina, F. M.  
1324 Maggi, Intra and inter-case features in predictive process monitoring: A  
1325 tale of two dimensions, in: *International Conference on Business Process*  
1326 *Management*, Springer, pp. 306–323.
- 1327 [33] K. Winter, F. Stertz, S. Rinderle-Ma, Discovering instance and process  
1328 spanning constraints from process execution logs, *Information Systems*  
1329 89 (2020) 101484.
- 1330 [34] W. Fdhila, M. Gall, S. Rinderle-Ma, J. Mangler, C. Indiono, Classifica-  
1331 tion and formalization of instance-spanning constraints in process-driven

- 1332 applications, in: International Conference on Business Process Manage-  
1333 ment, Springer, pp. 348–364.
- 1334 [35] C. W. Günther, H. M. W. Verbeek, XES standard definition, Technical  
1335 report, Eindhoven University of Technology, Eindhoven, The Nether-  
1336 lands, 2014.
- 1337 [36] T. Baier, J. Mendling, M. Weske, Bridging abstraction layers in process  
1338 mining, *Information Systems* 46 (2014) 123–139.
- 1339 [37] J. De Weerdt, M. De Backer, J. Vanthienen, B. Baesens, A multi-  
1340 dimensional quality assessment of state-of-the-art process discovery al-  
1341 gorithms using real-life event logs, *Information Systems* 37 (2012) 654–  
1342 676.
- 1343 [38] A. Augusto, R. Conforti, M. Dumas, M. La Rosa, F. M. Maggi, A. Mar-  
1344 rrella, M. Mecella, A. Soo, Automated discovery of process models from  
1345 event logs: review and benchmark, *IEEE Transactions on Knowledge  
1346 and Data Engineering* 31 (2018) 686–705.
- 1347 [39] M. J. Zaki, Spade: an efficient algorithm for mining frequent sequences,  
1348 *Machine learning* 42 (2001) 31–60.
- 1349 [40] H. Wickham, G. Grolemund, R for data science: import, tidy, transform,  
1350 visualize, and model data, O’Reilly, Sebastopol, 2017.
- 1351 [41] G. Janssenswillen, B. Depaire, M. Swennen, M. Jans, K. Vanhoof, bu-  
1352 par: enabling reproducible business process analysis, *Knowledge-Based  
1353 Systems* 163 (2019) 927–930.
- 1354 [42] L. Pufahl, T. Y. Wong, M. Weske, Design of an extensible BPMN  
1355 process simulator, *Lecture Notes in Business Information Processing*  
1356 308 (2017) 782–795.
- 1357 [43] F. Mannhardt, M. de Leoni, H. A. Reijers, W. M. P. van der Aalst,  
1358 P. J. Toussaint, From low-level events to activities - a pattern-based  
1359 approach, *Lecture Notes in Computer Science* 9850 (2016) 125–141.
- 1360 [44] M. Song, W. M. P. van der Aalst, Supporting process mining by showing  
1361 events at a glance, in: *Proceedings of the 17th Annual Workshop on  
1362 Information Technologies and Systems*, pp. 139–145.

- 1363 [45] I. A. Myrstad, Prosessstøtte i sengetun, Master's thesis, NTNU, 2017.
- 1364 [46] M. R. Halvorsen, H. O. Austad, A. D. Landmark, D. Ausen, I. Svagård,  
1365 T. Tomasevic, T. Trondsen, Redesigning work with a lightweight ap-  
1366 proach to coordination technology, CIN: Computers, Informatics, Nurs-  
1367 ing 37 (2019) 124–132.