

# Prediction of business process durations using non-Markovian stochastic Petri nets

Andreas Rogge-Solti<sup>a,\*</sup>, Mathias Weske<sup>b</sup>

<sup>a</sup>*Vienna University of Economics and Business, Austria*

<sup>b</sup>*Hasso Plattner Institute, University of Potsdam, Germany*

---

## Abstract

Companies need to efficiently manage their business processes to deliver products and services in time. Therefore, they monitor the progress of individual cases to be able to timely detect undesired deviations and to react accordingly. For example, companies can decide to speed up process execution by raising alerts or by using additional resources, which increases the chance that a certain deadline or service level agreement can be met. Central to such process control is accurate prediction of the remaining time of a case and the estimation of the risk of missing a deadline.

To achieve this goal, we use a specific kind of stochastic Petri nets that can capture arbitrary duration distributions. Thereby, we are able to achieve higher prediction accuracy than related approaches. Further, we evaluate the approach in comparison to state of the art approaches and show the potential of exploiting a so far untapped source of information: the *elapsed time* since the last observed event. Real-world case studies in the financial and logistics domain serve to illustrate and evaluate the approach presented.

*Keywords:* business processes, duration prediction, risk control, stochastic Petri nets

---

---

\*Corresponding author

*Email addresses:* andreas.rogge-solti@wu.ac.at (Andreas Rogge-Solti),  
mathias.weske@hpi.de (Mathias Weske)

## 1. Introduction

Managing a company's business processes is critical for its success in a competitive market environment [1; 2]. Not only is it important to strategically align the business processes with the company's goals, the daily operational support of business processes is also necessary to ensure smooth operation [3].

Psychological studies show that waiting time has a negative impact on customer satisfaction [4]. While companies cannot completely avoid waiting times for economic considerations, they can and should use all available information to provide their customers with reliable estimates of remaining time. An approach is to offer waiting time guarantees to customers. Kumar et al. show that these guarantees improve customer satisfaction, if they are met [5]. Such waiting time guarantees are common in various domains (e.g., fast food, travel, finance), where customers get compensated if their waiting time exceeds the agreed thresholds.

The increasing degree of automation by process oriented information systems and the adoption of sensing devices in business processes [6] produce large amounts of process execution data. These data hold valuable information which can be used for analyses in the context of business process intelligence. This paper focuses on the temporal performance aspect of business processes. In fact, we use expressive probabilistic models that we enrich with information extracted from event data [7]. Most prominently, performance models are used to predict remaining durations [8; 9; 10] and to estimate the risk of missing a given deadline [11; 12]

Based on these insights, we offer means to accurately predict remaining durations, and to compute risks to breach temporal deadlines or guarantees. Thereby, we exploit the information of elapsed time since the last event to increase accuracy. These means can be also used to compute reasonable waiting time guarantees, e.g., the waiting time that is met in 99 percent of the cases. Moreover, the approach presented in this paper can be applied in resource management, where an accurate predictor and knowledge about the uncertainty for remaining time of activities is critical for effective scheduling [13].

Summarizing the contributions of this paper, we rely on an expressive stochastic performance model of a business process to:

- predict the expected remaining duration of a business process,

- 38 • predict the risk of breaching a temporal deadline,
- 39 • use elapsed time as a means to improve prediction accuracy,
- 40 • evaluate the model against state of the art approaches.

41  
42 The remainder of the paper is structured as follows. In Section 2, we formally introduce the concepts and the model used in the prediction approach with a motivating example. Next, in Section 3, we present the approach to predict remaining time including associated confidence intervals, and also predict risk of breaching a temporal boundary. We overview related approaches in Section 4 and highlight differences to and synergies with our approach. Subsequently, Section 5 provides an evaluation with the state of the art and discusses the results. We conclude the paper in Section 6 with reflections on limitations and next steps.

## 51 **2. Preliminaries and Motivating Example**

52 To best illustrate the approach for prediction, we introduce a running example from the finance domain. Figure 1 depicts a loan approval process as a BPMN [14] diagram. In this process, there are two communicating partners: the client and the financial department of a bank. The client can start a loan process by submitting a loan request to the financial department. After that, the client waits for a reply. Meanwhile, the request is evaluated in the financial department and, depending on the evaluation, the loan is either accepted or declined. If declined, the client is notified and the process terminates. If the loan request is accepted, it first must be finalized and then can be approved by sending a notification.

62 In this process, the financial department has implemented a service level agreement (SLA) that states that the response will be sent to the client latest after 5 business days. In case of a later response the client is offered a compensation. Also, the management would like to have an estimate of the remaining case durations without having to ask the process participants to avoid excess communication. The vision is to also make the remaining duration estimates available to the affected clients of the bank, to further increase customer satisfaction. In this context, we investigate two issues: Predicting the remaining duration, and estimating the risk of deadline transgression in business processes.

72 Next, we define event logs that reflect the data gathered from information systems or sensors in a process execution environment. Let  $A$  be a set of

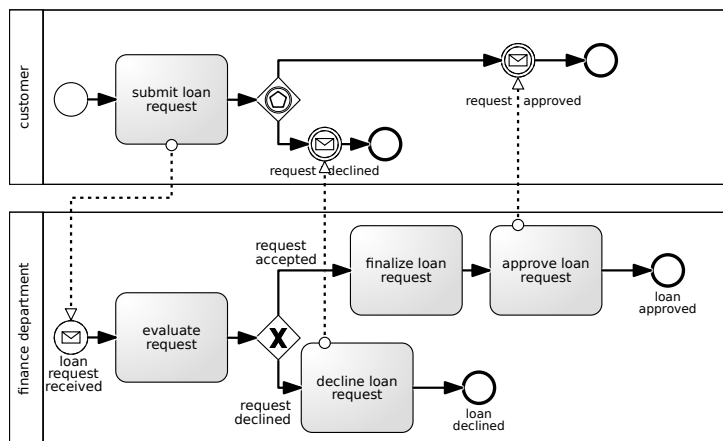


Figure 1: Loan application process as BPMN model. After a client requests a loan by submitting a request to the finance department, it is first evaluated. Depending on the evaluation, the loan can either be declined, or accepted. If accepted, the loan request will be finalized and an approved. The client is notified in any case about the decision.

74 process activities and  $TD$  be the time domain.

75 **Definition 1 (Event Log).** An event log over a set of activities  $A$  and time  
 76 domain  $TD$  is defined as  $L_{A,TD} = (E, C, \alpha, \beta, \gamma)$ , where:

- 77 •  $E$  is a finite set of events
- 78 •  $C$  is a finite set of cases (process instances),
- 79 •  $\alpha : E \rightarrow A$  is a function assigning each event to an activity.
- 80 •  $\beta : E \rightarrow C$  is a surjective function assigning each event to a case.
- 81 •  $\gamma : E \rightarrow TD$  is a function assigning each event to a timestamp. Events  
 82 are ordered by their timestamps.

83 We call the events belonging to the same case a *trace*. That is, the trace  
 84 of case  $c$  is  $\{e \in E \mid \beta(e) = c\}$ .

85 Beside the data, we need a model that probabilistically captures dura-  
 86 tions. The latter will be used to predict remaining durations and the risk of  
 87 deadline transgression. There exist many modeling notations (e.g., Business  
 88 Process Model and Notation [14]) that imperatively describe the permissible  
 89 behavior of business processes. We rely on the Petri net [15] representation  
 90 of these models, because Petri nets provide a formal model with a mathemat-  
 91 ical flavor and are able to capture most common control flow patterns [16].  
 92 That is, they are able to capture sequential dependencies, exclusiveness, and

93 concurrency. Moreover, there exist translations from different modeling lan-  
94 guages to Petri nets [17], which allow us to effectively apply our prediction  
95 approach in various settings. Hence, we can apply the approach independent  
96 of the modeling language that is used to capture the business processes by a  
97 prior translation of that model into a Petri net.

98 We first revisit the basic Petri net model, as introduced by Petri in his  
99 seminal thesis [15].

100 **Definition 2 (Petri Net).** *A Petri net is a tuple  $PN = (P, Tr, F, M_0)$*   
101 *where:*

- 102 •  *$P$  is a set of places.*
- 103 •  *$Tr$  is a set of transitions.*
- 104 •  *$F \subseteq (P \times Tr) \cup (Tr \times P)$  is the flow relation.*
- 105 •  *$M_0 \in P \rightarrow \mathbb{N}_0$  is an initial marking.*

106 Petri nets encode structural relationships between transitions, e.g., the  
107 firing of a transition creates tokens on its output places, which enable further  
108 transitions. To capture workflow systems, a desirable property of Petri nets  
109 was proposed by van der Aalst [16]. Therefore, he introduced the notion of  
110 workflow nets. A workflow net has two special places: a *source place*, and  
111 a *sink place*. The definition requires that all transitions and places of the  
112 net lie on a path between these two places. A token on the source place  
113 signifies the start of a case, and when a token reaches the sink place the case  
114 is completed. A workflow net is called *sound* [18], if any case will eventually  
115 terminate and the moment it reaches completion (with a token on the sink  
116 place) all other places are empty. Additionally, there are no dead transitions,  
117 i.e., for each transition, there is a firing sequence that includes it. In this  
118 paper, we assume that the process models are sound workflow nets.

119 Besides these properties, we need to capture performance characteristics  
120 of individual cases. For that purpose, we enrich transitions with proba-  
121 bilistic delays to reflect activity durations and waiting times. Moreover, we  
122 assign weights to transitions to capture the probability of decision outcomes.  
123 For example, to accurately predict the remaining duration, it helps to know  
124 that a path that usually takes very long, is only taken by 10 percent of the  
125 cases, while the other 90 percent follow the regular faster process path. The  
126 enriched model for our purposes is called generally distributed transition  
127 stochastic Petri net (GDT-SPN) [7], which we define as follows.

128 **Definition 3 (Generally Distributed Transition Stochastic Petri Net).**

129 *A GDT.SPN is a seven-tuple:  $\text{GDT.SPN} = (P, Tr, \mathcal{P}, \mathcal{W}, F, M_0, \mathcal{D})$ , where*  
130  *$(P, Tr, F, M_0)$  is the basic underlying Petri net. Additionally:*

- 131 • *The set of transitions  $Tr$  is partitioned into immediate transitions  $Tr_i$*   
132 *and timed transitions  $Tr_t$ .*
- 133 •  *$\mathcal{P} : Tr \rightarrow \mathbb{N}_0$  is an assignment of priorities to transitions, where  $\forall tr \in$*   
134  *$Tr_i : \mathcal{P}(tr) \geq 1$  and  $\forall tr \in Tr_t : \mathcal{P}(tr) = 0$ .*
- 135 •  *$\mathcal{W} : Tr_i \rightarrow \mathbb{R}^+$  assigns probabilistic weights to the immediate transi-*  
136 *tions.*
- 137 •  *$\mathcal{D} : Tr_t \rightarrow D$  is an assignment of arbitrary probability distribution*  
138 *functions  $D$  to timed transitions, reflecting the durations of the corre-*  
139 *sponding activities.*

140 Similar stochastically enriched versions of Petri nets are known in lit-  
141 erature, e.g., stochastic Petri net (SPN) [19], generalized stochastic Petri  
142 net (GSPN) [20]. The most basic form of probabilistic models, i.e., an  
143 SPN [19; 21; 22], is restricted to the use of only exponentially distributed  
144 durations. SPN models are efficiently analyzable due to being isomorphic to  
145 Markov Chains, which can be efficiently analyzed due to their *memoryless*  
146 property. That is, the exponential nature of transition durations abstracts  
147 from the fact how long a case has been waiting in a certain state, and there  
148 is only a rate of occurrence of the next event, which is constant. GSPN mod-  
149 els allow to model immediate transitions with weights, which do not have a  
150 stochastic delay associated, but fire immediately. GSPN mdoels have also  
151 been shown to be isomorphic to Markov chains [20] and can be analyzed  
152 efficiently.

153 In contrast, we allow arbitrary distributions to better model durations  
154 which need at least a minimum amount of time, or to capture deterministic  
155 timeouts. GDT.SPN models correspond to what Horváth et al. call non-  
156 Markovian stochastic Petri nets [23] in their analysis.

157 We have shown that we can enrich Petri nets with stochastic execution  
158 information gathered from data [24]. In that work an implementation using a  
159 maximum likelihood approach is provided open-source in the process mining  
160 platform ProM [25] in the package `StochasticPetriNets`. The idea is to fit  
161 distributions to the observed durations at each transition in the model and  
162 learn also the probabilities of decisions in the model.

163 Now we return to the motivating example. We assume that the financial  
164 department has created a business process model for the loan application

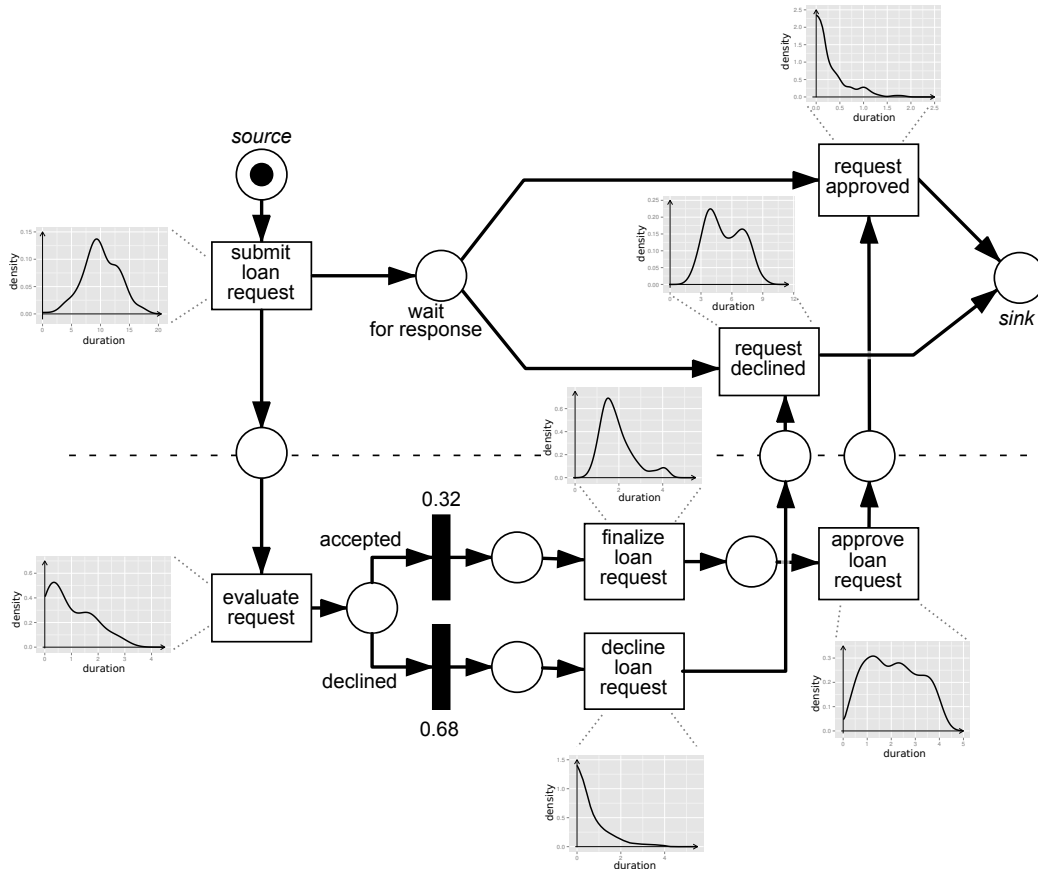


Figure 2: Loan application process as GDT.SPN model. Converted from the BPMN model in Figure 1, cf. [17] and enriched with performance information [24]. Places are circles and transitions are boxes. The model has a *source* and *sink* place indicating start and end. Each timed transition is a white box and is annotated with the corresponding duration distributions. The immediate transitions are filled black bars and are annotated with their weights (e.g.,  $\mathcal{W}(t_{accept}) = 0.32$ ) that reflect the probability of firing. That is, loan requests are accepted in 32 percent of the cases, whereas 68 percent get declined. The horizontal dashed line has no semantical meaning, it highlights the boundary of communication between client (above) and financial department (below). The probability density functions associated to timed transitions represent the distribution of firing durations.

165 process (see Figure 1). In this case, we can either manually, or with the help  
 166 of automatic translation, convert the BPMN process model to a Petri net  
 167 model [17]. If we collect and replay historical process execution data that can  
 168 be extracted from the supporting IT systems, as presented in [24], we obtain a  
 169 GDT.SPN model. The output model for the running loan application process

170 example is depicted in Figure 2.

171 Given a GDT.SPN model that contains decision probabilities and transi-  
172 tion durations, and given a running case of the process (i.e., the events that  
173 have been produced so far by supporting IT systems), we can compute the  
174 remaining case durations and risks of deadline transgressions. In the next  
175 section, we shall discuss the core of this contribution: the conceptual advan-  
176 tage of using *elapsed time* since the last event, and the algorithm that allows  
177 us to make estimates about the remaining time in GDT.SPN models.

### 178 3. Prediction considering elapsed time

179 In this section, we present the conceptual solution to predicting remaining  
180 durations and estimating risk of deadline transgression with a framework  
181 based on simulation. The estimation of the remaining time to reach a given  
182 state is formalized as follows.

#### 183 3.1. Estimating the remaining time

184 Let  $t_M \in TD$  be the time when the process reaches the target marking  $M$   
185 of interest (e.g., the final marking). We do not know  $t_M$  in advance, because  
186 it is subject to random effects. Hence, we capture it by a random variable  
187  $X_M$ . Let  $t_0$  be the current time at estimation, and  $t_a$  be the time of finishing  
188 the current activity  $a$  for a given case. Then, estimating remaining time is to  
189 compute the expected value of the remaining relative duration  $E(X_M - t_0)$ .  
190 In [10], we analyzed the accuracy gained by conditioning on the durations of  
191 the current activity in the process. Concretely, we compared the two esti-  
192 mates  $E(X_M - t_0 | t_a \geq t_0)$ , and  $E(X_M - t_0)$ . Note that if there are multiple  
193 activities concurrently running, their durations are conditioned accordingly.  
194 For example, let  $t_b$  be the duration of a second activity  $b$  running at time  $t_0$ .  
195 Then, we compute  $E(X_M - t_0 | t_a \geq t_0, t_b \geq t_0)$

196 The resulting estimations can be used in various ways. For example,  
197 patients can be informed about their expected length of stay, or the clients  
198 awaiting response can be notified with the expected waiting duration.

#### 199 3.2. Estimating the risk of transgression

200 This section is about the probability of meeting a certain deadline. There-  
201 fore, let  $t_d$  be the deadline that should be met. Let again  $t_M \in TD$  be the  
202 time of reaching the desired marking. Thus, we are interested in estimating  
203 the probability  $P(t_M > t_d)$  that reflects transgressing the deadline. For this



204 problem, we also can condition the duration of the current activity on the  
 205 elapsed time. This results in the formula  $P(t_M > t_d \mid t \geq t_0)$ , cf. Section 3.1.  
 206 Note that while the first question yields a relative duration (which is a value  
 207 greater or equal to 0), the second formula yields a probability between the  
 208 two extremes of 0 (impossible) and 1 (always).

209 Based on these estimations, a process controller can find cases in risk of  
 210 transgressing a deadline, and decide whether to take action to speed up exe-  
 211 cution. Based on how risk-averse the management is, and the costs inflicted  
 212 by a transgression of the deadline, the process controller can decide to control  
 213 the process when a certain risk threshold is exceeded. For example, it can be  
 214 optimal in terms of cost to only act, when the risk of exceeding a deadline  
 215 becomes higher than 60 percent.

### 216 3.3. Conceptual information gain by capturing elapsed time

217 When using a monitoring system that records process steps in an au-  
 218 tomated way, the events recorded are instantly available to the monitoring  
 219 system. Our assumption in this work is that the durations of process steps  
 220 (e.g., human tasks, web service calls) is much greater than the time taken to  
 221 register the event in the system.

222 Based on this assumption, a previously neglected source of information  
 223 can be tapped and exploited to get a more accurate prediction model [10]:  
 224 Elapsed time since the last event. We illustrate the effects of conditioning  
 225 on remaining time with characteristic example duration distributions.

226 Formally, let  $t \in TD$  represent time. Let  $tr \in Tr_t$  be a timed transition  
 227 with the assigned duration distribution function  $F_\delta = \mathcal{D}(tr)$ . We obtain the  
 228 *density* function  $f_\delta$  by differentiating the distribution function  $F_\delta$ , that is,  
 229  $f_\delta(t) = dF_\delta(t)/dt$ . Let  $t_0 \in TD, t_0 \geq 0$  be the current time since enabling  
 230 of  $tr_i$ . Let further  $f_{\delta_{Dirac}}$  denote the Dirac delta function, which captures the  
 231 whole probability mass at a single point. We define the density function of  
 232 the truncated distribution as:

$$f_\delta(t \mid t \geq t_0) = \begin{cases} 0 & t < t_0 \quad F_\delta(t_0) < 1 \\ \frac{f_\delta(t)}{1-F_\delta(t_0)} & t \geq t_0, \quad F_\delta(t_0) < 1 \\ f_{\delta_{Dirac}}(t - t_0) & F_\delta(t_0) = 1 \end{cases} \quad (1)$$

233 The part of the density function that is above the threshold  $t_0$  is rescaled  
 234 such that it integrates to 1, which is a requirement for probability density  
 235 functions. Note that in the exceptional case that  $F_\delta(t_0) = 1$  (i.e., the current

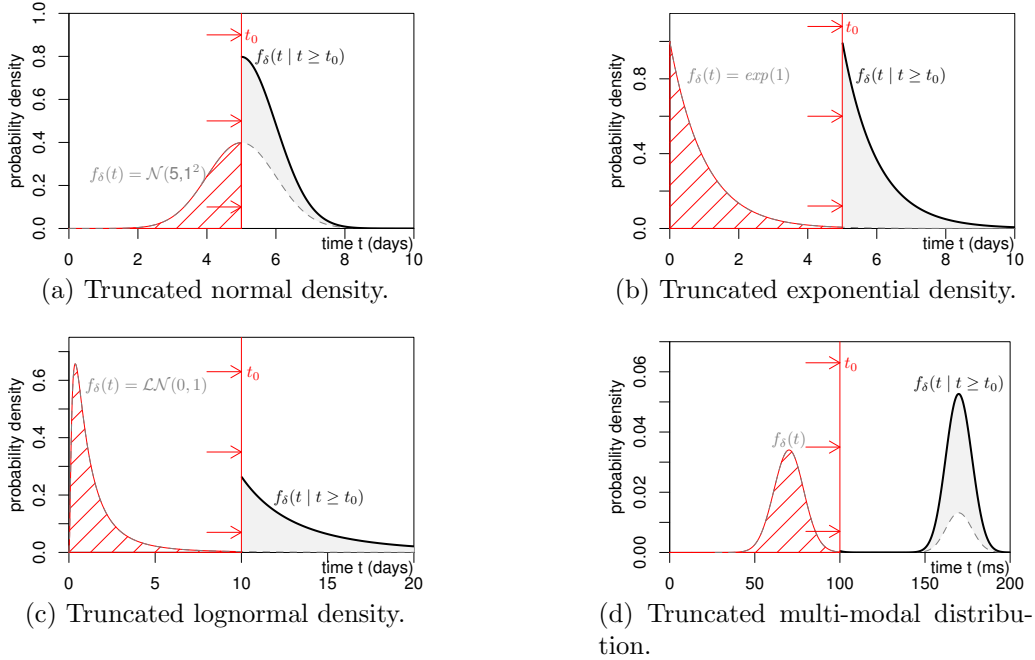


Figure 3: Different truncated probability density functions [7]. The current time  $t_0$  truncates the distributions and is depicted as a vertical line.

236 time  $t_0$  progressed further than the probability density function’s support),  
 237 we use the Dirac delta function with its peak at  $t_0$ . In this case, the activity  
 238 is expected to finish immediately at  $t_0$ .

239 The intuition is as follows. We base our predictions on a stochastic model  
 240 describing the distribution of a large amount of cases. Thereby, we discard  
 241 the fraction of the cases that is not consistent with our observation for the  
 242 current activity’s duration, i.e., those cases that would have completed the  
 243 currently running activity before the current time. In contrast to using condi-  
 244 tional probability density functions, traditional methods predict the remain-  
 245 ing duration of a case only at event arrival, and subtract elapsed time from  
 246 the predicted duration at later points in time [26; 8].

247 Figure 3 shows the effect that *truncation* has on different types of dis-  
 248 tributions. The density of the normal distribution (depicted in Figure 3a)  
 249 decreases faster than that of the exponential distribution. Therefore, trunca-  
 250 tion of the normal distribution makes us more certain that the corresponding  
 251 event will happen soon. In the figure, the value of the conditional density

252 function is higher than the original density and more concentrated.

253 The exponential distribution, depicted in Figure 3b, is a special case,  
254 where conditioning on elapsed time does not affect the shape of the distri-  
255 bution, that is,  $exp(t) = exp(t + t_0 | t > t_0)$ . Latter property is an effect of  
256 the memoryless property. It implies that the probability that an event will  
257 occur in the next minutes is not affected by the time that we spent waiting  
258 before. That is, if  $X$  is an exponentially distributed random variable, then  
259 for any  $t, t_0 \geq 0$  it holds that  $P(X > t_0 + t | t > t_0) = P(X > t)$ .

260 Heavy-tailed distributions (e.g., the lognormal distribution depicted in  
261 Figure 3c) are on the other side of the spectrum. Latter distributions can get  
262 a higher variance, if we condition on elapsed time. This means that the longer  
263 we wait for an event, the less likely it will be that we will observe it in the next  
264 minutes. The uncertainty and the expected value of the distribution grows  
265 with the time that we spend waiting. In the figure the value of the conditional  
266 density is lower and decreasing at a slower rate than the unconditional one.

267 We can gain the most information about the expected events, when condi-  
268 tioning multi-modal distributions (i.e., distributions with multiple peaks in  
269 their density functions). An example is depicted in Figure 3d. The depicted  
270 distribution has two modes, that is, its values are clustered at two different  
271 time points. Here, once the time progressed beyond the first mode, it can  
272 be assumed that the duration will be around 170 milliseconds. This example  
273 shows that we can get more accurate results by excluding inconsistent cases  
274 from prediction, especially if we use non-parametric methods that are able  
275 to approximate distributions with multiple modes.

### 276 3.4. Prediction algorithm

277 Besides the already mentioned assumption of immediate detection of  
278 events by the prediction framework, we consider each activity duration in  
279 isolation, independently from other activities. This is a common simplifying  
280 assumption that we share with all analytical approaches to prediction.

281 To make predictions for a single case, we require the current state of the  
282 case, and the model that captures experiences about the behavior of the  
283 process. The prediction algorithm takes five inputs:

- 284 1. the GDT.SPN *model* of the business process, cf. Definition 3,
- 285 2. the current time  $t_0$ ,
- 286 3. the deadline  $t_d$ ,
- 287 4. the current *trace* of the case  $c$ . That is, all observed events up to  
288 time  $t_0$ . More explicitly, this is  $\{e \in E \mid \beta(e) = c \wedge \gamma(e) \leq t_0\}$ ,

289 5. the number of simulation iterations  $n$  controlling the precision of the  
 290 prediction.

---

**Algorithm 1** Prediction algorithm

---

```

1: procedure PREDICT(model,  $t_0$ ,  $t_d$ , trace,  $n$ )
2:   currentMarking  $\leftarrow$  replay(trace, model)            $\triangleright$  find current state
3:   times  $\leftarrow$  new List()                                $\triangleright$  used to collect results
4:   for all  $i \in \{1, \dots, n\}$  do
5:     time  $\leftarrow$  simulateConstrained(model, currentMarking,  $t_0$ )
6:     times.add(time)
7:   end for
8:   return getMean(times), getRatioAbove(times,  $t_d$ )
9: end procedure

```

---

291 Algorithm 1 describes the procedure. It starts by finding the appropriate  
 292 current state, i.e., the marking, in the model by replaying the available ob-  
 293 served events of the case in the model. We assume that a workflow engine is  
 294 in charge of controlling the process flow, which facilitates replay of observed  
 295 events in the log. If this is not the case, we support an alignment tech-  
 296 nique [27] for replaying non-fitting traces on the model to find the current  
 297 state.

298 Next, the algorithm collects simulation results, i.e., completion *times*, of  
 299 the given number of simulation iterations  $n$  in a list. Each simulation run  
 300 represents a sample from the possible continuations of the process according  
 301 to the model. Each individual sample is independent and identically dis-  
 302 tributed, due to the independence assumptions in the model, cf. Definition 3.  
 303 The **simulateConstrained** method simulates continuations of the trace for  
 304 the GDT-SPN model, but instead of sampling from the original transition  
 305 distributions  $F_\delta(t)$ , it samples from the truncated distributions conditioned  
 306 on the current time  $F_\delta(t \mid t \geq t_0)$ , as described in Section 3.3.

307 The maximum likelihood estimate for the remaining time is the sample  
 308 mean of the duration samples stored in the list *times* with size  $n$ . For conve-  
 309 nience, we can access each element in the list with an index  $i \leq n$ , e.g., *times*<sub>3</sub>  
 310 is the third duration sample. Then, we approximate  $E(t_X - t_0 \mid t \geq t_0)$  by  
 311 the sample duration values, i.e., we compute:

$$\widehat{E}(t_X - t_0 \mid t_{next} \geq t_0) = \frac{1}{n} \sum_{i=1}^n [times_i - t_0] \quad (2)$$

312 The above measure is a point estimate to the remaining mean duration  
 313 that does not include our degree of certainty in the measure. We can also  
 314 compute the confidence interval estimate for the remaining mean duration  
 315 to be able to say that in most cases (e.g., in 95 percent of the estimations)  
 316 the interval will contain the real remaining mean duration [28, Chapter 6].

317 Note that the accuracy of a prediction based on simulated samples de-  
 318 pends on both the number of computed samples, as well as the standard  
 319 deviation within the samples. Therefore, we also support the mode, where  
 320 instead of a sample size, the user can set required accuracy thresholds. For  
 321 example, the simulation continues taking samples, until the 99 percent con-  
 322 fidence interval on the prediction lies within  $\pm 3$  percent of the predicted  
 323 value.

324 For predicting the chance to transgress a target deadline  $t_d$ , we use the  
 325 maximum likelihood predictor. This is simply the ratio of the samples taking  
 326 more time than the deadline to the number of samples. Let *times* again be  
 327 our list of samples with size  $n$ . Then, we estimate the probability to miss a  
 328 given deadline  $t_d$  as follows:

$$\widehat{P}(t_X > t_d) = \frac{|\{t \in \text{times} \mid t > t_d\}| + \epsilon}{n + 2\epsilon} \quad (3)$$

329 Note that maximum likelihood methods are prone to overfitting the sam-  
 330 ple, e.g., it is usually not advised to have extreme probabilities of 0 to out-  
 331 comes that are possible, but were not observed so far. Therefore, we apply  
 332 *Laplace smoothing* for the estimator with the parameter  $\epsilon = 0.5$  which avoids  
 333 extreme probability values. Chen and Goodman provide an overview of com-  
 334 mon smoothing techniques [29].

### 335 3.5. Open-Source Implementation

336 We implemented the prediction algorithm in the process mining frame-  
 337 work ProM as a plugin<sup>1</sup>. The method to enrich a Petri net with historical  
 338 performance data extracted from a log is also available in that plugin [24].  
 339 It is possible to learn different kind of parametric models (e.g., normally  
 340 distributed durations), as well as nonparametric models (e.g., simple his-  
 341 tograms, or kernel density estimators). If the learned models are only used

---

<sup>1</sup>Implementation provided open-source in the `StochasticPetriNet` package of ProM.  
 Available at <http://www.promtools.org>

342 for prediction, simple histograms based on the observed samples suffice for  
343 making predictions. In this latter case of histograms, the sampling method  
344 can simply pick one of the past observations randomly. We exclude the ob-  
345 servations that do not meet the constraint of being greater or equal to the  
346 current time  $t_0$ . There might be cases, when the current instance takes longer  
347 for an activity than all previously observed cases. In this cases, the histogram  
348 based sampling returns the constraint  $t_0$  itself.

349 When expert estimates exist and parametric probability distributions are  
350 used in the GDT.SPN model, e.g., normal, exponential, or lognormal distri-  
351 butions, we use rejection sampling. Rejection sampling is simple: we draw a  
352 sample from the original distribution and check, whether it meets the given  
353 constraints. If not, we reject the sample and repeat the process until we get  
354 a sample that meets the constraints.

355 Before evaluating the prediction method with real data, we introduce  
356 related approaches against which we will benchmark the prediction quality.

#### 357 **4. Related Approaches**

358 There are many dimensions that one can predict (e.g., time, costs, re-  
359 source allocation), however, we focus on the time aspect in this paper. Con-  
360 cretely, we want to predict (i) the remaining time of a current case until it  
361 reaches a target state, and (ii) the risk of missing a given temporal deadline  
362 for a target state. Therefore, we group related approaches to either of these  
363 two categories. First, we discuss approaches predicting remaining time.

##### 364 *Remaining Time Prediction Approaches*

365 In their work, van der Aalst et al. [8] use the available information in logs  
366 to predict the remaining case duration based on observed durations in the  
367 past. They first create an annotated transition system from the event log. In  
368 that state transition system, they collect remaining durations for each visited  
369 state from the traces in the log. Finally, when a given case is executed, in  
370 each state the remaining time is predicted as the average remaining duration  
371 of the historical cases that passed that state. Our approach is similar in  
372 the sense that it also abstracts from data and resources, but uses GDT.SPN  
373 models instead of transition systems, making our approach more accurate  
374 when parallelism exists in the process. The work in [8] serves as one of the  
375 benchmarks for the prediction method proposed in this paper.

376 Building on the work in [8], Folino et al. [9] present an improvement based  
377 on predictive clustering. They make use of additional contextual information  
378 of a trace (e.g., the current workload in the system) to perform clustering.  
379 The idea is to group similar traces and base predictions for new ones with  
380 similar features on only similar historical cases of the log. They use the  
381 predictions to warn in case of a predicted transgression of a threshold. It  
382 seems promising to combine the work in [9] with a GDT\_SPN approach.

383 Other work for prediction of performance was presented by Hwang et  
384 al. [30] and similarly Zheng et al. [31]. They use formulae to compute quality  
385 of service criteria, such as expected durations of compositions. Typically,  
386 these works assume the service compositions to be composed of building  
387 blocks, that is, of single-entry single-exit blocks, cf. the formal definition  
388 in the work by Kiepuszewski et al. in [32]. The methods proposed can be  
389 used for business processes, too. However, the block-structured assumption is  
390 lifted in this work, allowing for more complex models, and we target already  
391 running instances.

392 Closely related to our approach is the prediction method presented by  
393 Wombacher and Iacob [33]. In their work, the authors prepare event logs of  
394 unstructured processes for prediction of activity durations. They use mean  
395 duration of activities for predictions, but they do not use runtime information  
396 of elapsed time since the last event.

397 Also based on the assumption that activity durations are normally dis-  
398 tributed, the work by Anklesaria et al. estimates completion time for PERT  
399 networks [34]. PERT networks are used to model projects with required ac-  
400 tivities that are in dependency relations, e.g., an activity can only begin after  
401 two previous activities have been completed. Optionality and loops are not  
402 captured in such models, however, and the remaining project time is deter-  
403 mined by the longest path through the network. Anklesaria et al. investigate  
404 the effects of correlations between different paths to increase the accuracy of  
405 the predicted error by taking into account multiple paths.

406 Recently, Senderovich et al. proposed to mine queueing networks from  
407 logs to enable remaining time prediction [35]. Their work yields high accu-  
408 racy as resources are captured in their model, but it also requires explicit  
409 information about waiting time, i.e., the entry time and exit time of queues  
410 at activity stations. A more general approach is followed by de Leoni et  
411 al. [36], where a framework is presented that maps the prediction problem in  
412 business processes to a well-known data mining setting. That is, they map  
413 each event of an event log to a table with several associated attributes to

414 learn associations and regressors. Thereby, they allow to incorporate further  
415 data. However, they do not treat parallelism in their approach. The result-  
416 ing predictor is similar to the method mentioned in [8], but with more data  
417 on which to condition the remaining time.

418 Simulation has also been proposed and used for operational decision mak-  
419 ing by Rozinat et al. [37; 38]. The idea is to set up a simulation environment  
420 capturing the current situation and start a short-term simulation from this  
421 state with different simulation parameters. Their use of simulation is for  
422 operational decision support and is focused on the overall performance of  
423 business processes. In contrast, we use simulation to make a prediction for  
424 the current instance only and use the current elapsed time as additional input  
425 to the simulation which allows to improve single predictions.

426 Analysis of non-Markovian stochastic Petri nets (i.e, GDT-SPN models  
427 in this paper) has already been done before. Monte Carlo simulation is the  
428 preferred choice for analysis, e.g., in [39; 40]. However, previous work does not  
429 address remaining durations of single instances with conditional probability  
430 densities.

#### 431 *4.1. Approaches dealing with process related risks*

432 The interest in risk-aware business process management is increasing since  
433 one of the earliest works put emphasis on the topic [41]. A recent survey on  
434 the topic is provided by Suriadi et al. [42]. While the earlier works approach  
435 the problem from a qualitative view, recently quantitative approaches were  
436 also proposed. The works of Pika et al. [12], Goluch et al. [43], and Conforti  
437 et al. [11] are closely related to our approach.

438 Pika et al. [12] define indicators for the risk of deadline violations. They  
439 search for patterns, such as abnormal activity durations, and use that infor-  
440 mation for predicting whether a case will be late. Goluch et al. design a pro-  
441 totype for risk-aware simulation with the special focus of simulating threats  
442 to resources, e.g., power outage, or illness of a process participant [43]. In  
443 contrast, we focus on the risk of missing a temporal deadline. Further, Kang  
444 et al. [44] advocate business process monitoring in real time. Their approach  
445 is based on classifying historical traces in correct and incorrect traces by  
446 data mining techniques, such as support vector machines. Their goal is also  
447 to predict and classify current instances, e.g., if they are likely to exceed  
448 deadlines. Their approach only captures sequential processes, however, and  
449 timestamps of events are not considered in the prediction.



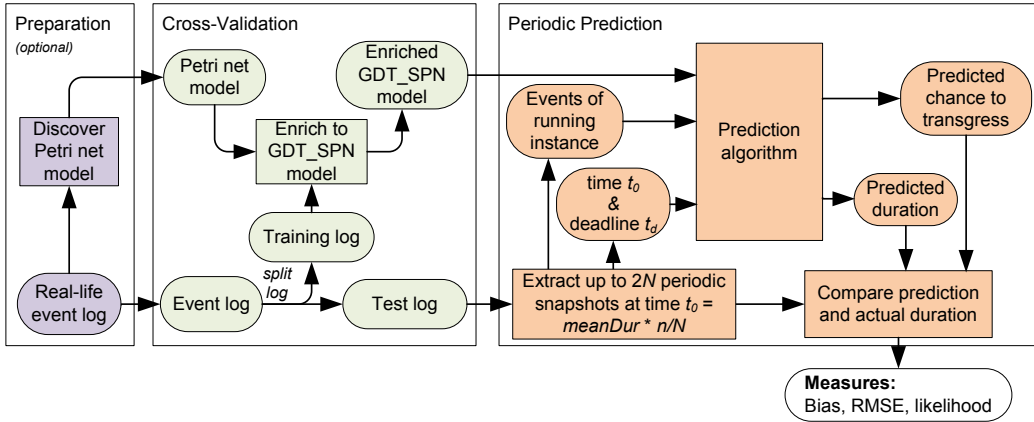


Figure 4: Evaluation setup. The figure contains an optional preparation phase, which we require if we do not have a Petri net model for real-life event logs. After preparation, i.e., simulation of a GDT\_SPN model and the extraction of its underlying Petri net, the cross validation is performed. Here, the log is split into training and test log. The training log is used to enrich the GDT\_SPN model. The test log is used for evaluating the predictions of the model. We measure prediction accuracy at  $2N$  equidistant points in time, where the mean case duration is at the  $N^{\text{th}}$  measurement.

450 The work presented by Leitner et al. [45] considers running instances.  
 451 They use regressions for durations between two-point measures in the process.  
 452 The predictions are then used to identify whether a service level agreement  
 453 will be violated. This method is By contrast, our work includes knowledge of  
 454 the whole business process model to make predictions and we use the elapsed  
 455 time since the last event as constraining factor.

456 In summary, the approach presented here is the first one to take elapsed  
 457 time into account for prediction of remaining durations. In the following, we  
 458 evaluate the approach with state-of-the art approaches.

## 459 5. Evaluation

460 The evaluation is performed in two parts: a qualitative analysis, and  
 461 a scalability analysis. In the qualitative analysis, we assess the prediction  
 462 quality of predicting the remaining durations of a process. We further assess  
 463 the estimation of the risk of transgressing a given deadline.

### 464 5.1. Experiment setup

465 The procedure works as depicted in Figure 4. Depending on whether we  
466 already have a Petri net model for an event log, we need a preparation phase  
467 to discover [46; 47] a suitable Petri net model representing the structure of  
468 a process.

469 Given the log and the Petri net input, we perform a 10-fold cross vali-  
470 dation with the event log, i.e., we split it into ten parts, use nine parts to  
471 enrich the Petri net model to a GDT-SPN model and use the left-out part as  
472 test log. The test log represents new cases as they would appear in reality.  
473 The experiment is repeated such that each of the ten parts are used once as  
474 test log.

475 For each trace in the test log, we perform  $2N$  periodic prediction itera-  
476 tions. Thereby, we run the  $N^{\text{th}}$  prediction *meanDur* after the start of the case  
477 (*meanDur* is the average case duration). In our experiments  $N$  is 20, such  
478 that we have 40 prediction iterations in total. Each time we perform a pre-  
479 diction, we supply the current time  $t_0$ , the target deadline  $t_d$ , and the current  
480 *trace* up to  $t_0$  to the algorithm. In our approach the model is a GDT-SPN  
481 model, but we also compare with the prediction approach proposed in [8].  
482 In latter case, the prediction model is a transition system mined from the  
483 event log. The approach in [8] allows different configurations of the transi-  
484 tion system. States can be very informative, e.g., encode the entire *sequence*  
485 of all events observed so far. But they can also be rather general, e.g., only  
486 encoding the single *last* observed event. Therefore, we compare to different  
487 configurations of the transition system based approach. As baseline, we also  
488 compare to a predictor that predicts the remaining time to be the difference  
489 to the mean case duration. In case a prediction would be negative (i.e., when  
490 the current time exceeded the mean duration) this predictor returns 0.

### 491 5.2. Loan approval process example

492 Our motivating example of Section 2 is inspired by the process that was  
493 presented in the Business Process Intelligence Challenge 2012 [48]. The loan  
494 application process of a Dutch financial institute was offered for public analy-  
495 sis. We depict the Petri net representation of the top level process in Figure 5.  
496 The process starts with a submission of a loan request. The request can be  
497 immediately declined, or be preaccepted by the system. Then, a clerk can  
498 accept the loan request to later finalize, and approve it. At these stages, the  
499 client also can cancel the loan request, or it can be declined by the clerk.

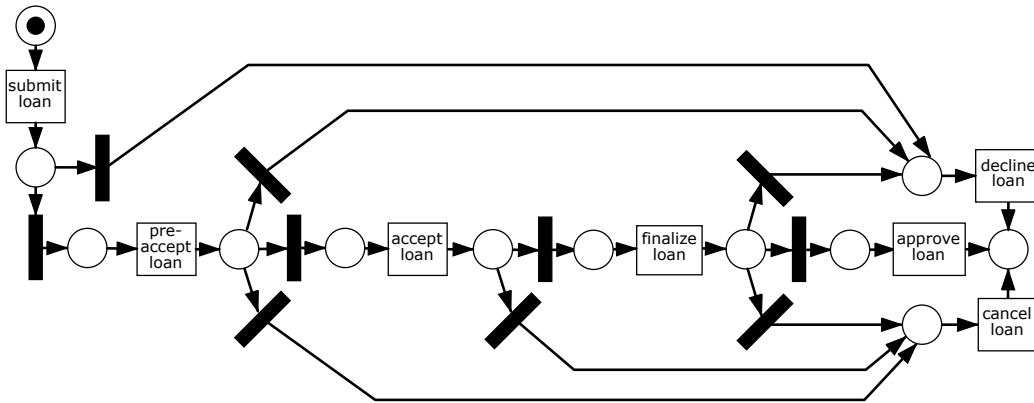


Figure 5: This Petri net process model captures a high level view on a loan application process. Besides the sequential *happy path* (i.e., submitting, preaccepting, accepting, finalizing, and approving the loan), there is also the option for the client to cancel the loan or for the clerk to decline it.

500 The average duration of all cases in this log is almost nine days, while many  
 501 cases are quickly declined by the system.

### 502 5.3. Shipment import process example

503 The second real-life process analyzed in our evaluation is from the logistics  
 504 domain. Figure 6 depicts the model of the shipment import process. Each  
 505 case reflects the events associated with one container that passes the Dutch  
 506 harbor to its destination in Europe. The event log contains event entries  
 507 for arrival by sea, discharge, and subsequent pick up of the containers at the  
 508 logistics provider. The average duration of the shipment process is four days.  
 509 Due to the high grade of automation in this process, the fitness between the  
 510 event log and the process model is perfect.

511 Giving an accurate estimation of the time of arrival of a container to  
 512 clients is helpful for their subsequent planning and allows to save warehousing

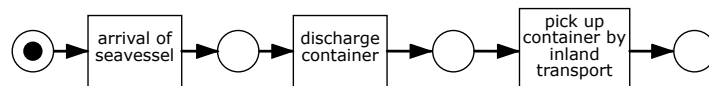


Figure 6: The process model captures three sequential stages of a logistics process: the arrival of sea vessels, container discharges, and further inland shipping.

513 costs by just in time delivery. It is equally important to have accurate models  
514 that can estimate the chance to miss a given shipment guarantee.

515 Note that the number of activities in this process are few (sparse moni-  
516 toring points) and therefore, conditioning on each activity’s duration should  
517 yield higher benefits compared to a process with a large number of activities  
518 (dense monitoring points).

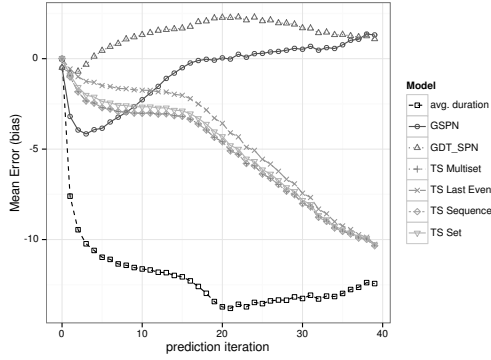
#### 519 5.4. Experiment results

520 To assess the proposed approach, we evaluate (i) the precision of the  
521 predictions, (ii) the accuracy of the predictions, and (iii) the model fit to  
522 predicting missing a given deadline. We measure each of these dimensions  
523 for the compared models in each of the 40 periodic predictions as explained  
524 in Section 5.1.

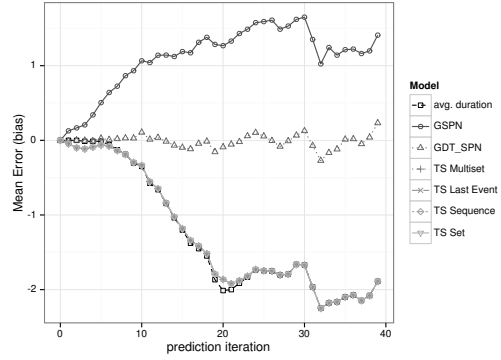
525 Precision (i) tells us how far the predictors are away from the actual  
526 value on average. We measure precision with the model *bias*, for which  
527 we use the expected error (i.e., the mean of the prediction errors) as an  
528 estimator. In other words, this measure tells us, if our model tends to over-  
529 or underestimate the remaining duration. A bias of 0 is optimal. A negative  
530 value indicates underestimation, a positive value indicates overestimation.

531 We measure the accuracy (ii) by the root mean square error (RMSE).  
532 The measure is the square root of the average squared errors. It is a common  
533 measure that tells us, how close the predicted value is to the real value. The  
534 lower this measure is, the more accurate is the prediction model.

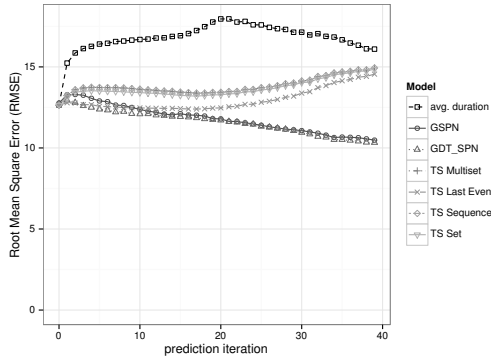
535 Last, we measure the model’s fit (iii) to predicting missing a given dead-  
536 line  $t_d$  by the *log-likelihood* of the model given the data. The deadline  $t_d$   
537 is twice the mean case duration for our experiments. In this setting, the data  
538 encode whether cases exceed the deadline. Models predict either of these  
539 exclusive outcomes with probabilities summing to 1. For example, a model  
540 can predict at a given point that the current case will exceed the deadline  
541 with a probability of 0.1, and be finished within the deadline with probability  
542 0.9. Let us assume that we observe 2 cases transgressing the deadline and  
543 3 cases meeting the deadline in this case, and that the probabilities remain  
544 the same. Then, the log likelihood is simply  $\log(0.1^2 \cdot 0.9^3) = -2.14$ . A  
545 better model in this example would predict that the chance to transgress  
546 is 0.4, and the chance to meet a deadline is 0.6. Latter model would yield  
547 a higher log likelihood of  $-1.46$ . The higher the log-likelihood, the better.  
548 Although the measure by itself is not too informative as it also depends on  
549 the number of observations, we can estimate a model’s relative performance



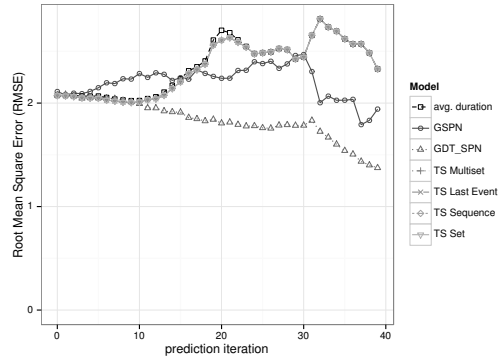
(a) Prediction *bias* in days for the loan approval process [48].



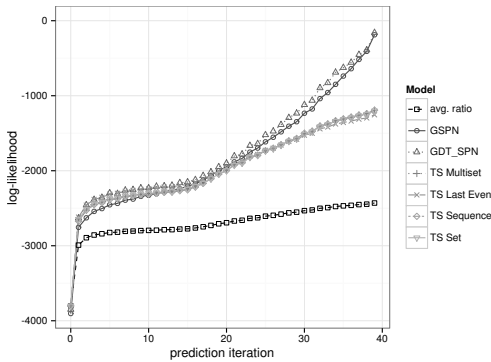
(b) Prediction *bias* in days for the shipment import process.



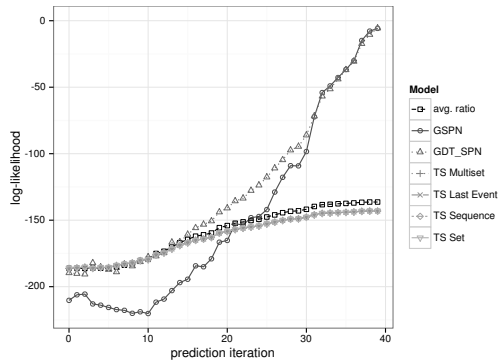
(c) Prediction *accuracy*. RMSE in days for the loan approval process.



(d) Prediction *accuracy*. RMSE in days for the shipment import process.



(e) Model log-likelihoods of deadline transgression for the loan approval process.



(f) Model log-likelihoods of deadline transgression for the shipment import process.

Figure 7: Prediction quality for the logistics process. The root mean square errors are collected for periodic prediction iterations of the case studies. Thereby, a 10-fold cross validation is performed on the data with 40 periodic predictions, s.t. the 20th iteration is at the mean duration. The deadline  $t_d$  is set to be twice the mean duration.

550 by comparing its log-likelihood with another model’s log-likelihood in each  
551 prediction iteration.

### 552 5.5. Discussion of the qualitative analysis

553 Based on the assumption that the time from occurrence of events to their  
554 detection is negligible in comparison to activity duration, we proposed to use  
555 the elapsed time as conditioning factor on the duration distributions. We  
556 compared various prediction algorithms on two real life event logs.

557 In Figure 7a, we see that the baseline model underestimates the dura-  
558 tions of the loan approval process. This can be explained by many cases  
559 automatically rejected at an early stage, but still considered in the remain-  
560 ing predictions. After a first underestimation, the GSPN model assuming  
561 exponential distributions performs very well in this case. Our GDT-SPN ap-  
562 proach based on truncated duration distributions tends to overestimate the  
563 remaining durations, but in later predictions yields comparable results to the  
564 GSPN model. The transition system based models start to underestimate  
565 the longer running cases around the 15<sup>th</sup> prediction iteration, as they do  
566 not properly condition on the remaining time. These effects are more pro-  
567 nounced in Figure 7b, where our approach is rather unbiased (i.e., the bias  
568 is close to 0). The GSPN model is not so suited, however, as the exponential  
569 distribution assumption overestimates remaining times in this process. The  
570 transition system based approaches cannot exploit the simple structure of  
571 the process, and collapse with the single last event predictor.

572 Regarding prediction *accuracy*, we see that in Figure 7c, the two Petri net  
573 based predictors start to outperform the transition system based predictors  
574 at later prediction iterations. The GSPN based predictor yields worse results  
575 in the shipment import process, cf. Figure 7d. Here, the implicit exponential  
576 assumption made by GSPN models does not hold. Note that our approach  
577 yields significantly smaller errors than the compared predictors, especially  
578 after some time has passed.

579 We can see in Figures 7e and 7f that the log likelihoods for all models  
580 increase over time. This is mostly caused by the decrease in the number  
581 of cases still running at later predictions. Nevertheless, the two Petri net  
582 approaches prove to have better fitness in providing answers to the question  
583 whether a case will take longer than twice the usual time. Especially at  
584 later prediction periods it is required to condition on elapsed time to predict  
585 whether the current case will take longer than the target deadline.

586

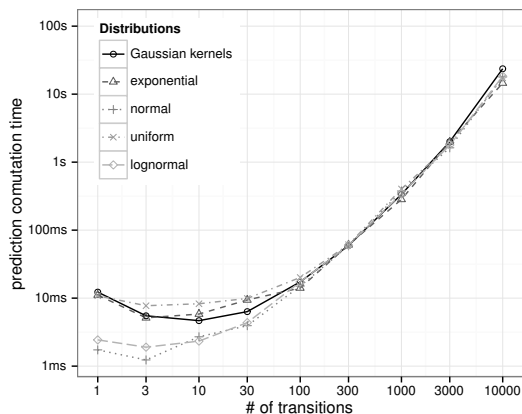


Figure 8: Prediction algorithm durations for differently sized randomly generated, acyclic GDT\_SPN models. Axes in log-scale.

587 To summarize the qualitative evaluation results, the following character-  
 588 istics of our approach can be observed:

- 589 • At early prediction periods, the proposed approach performs about as  
 590 well as the benchmarks.
- 591 • The improvements become more pronounced as time proceeds.
- 592 • The approach is most valuable for long running cases, which are critical  
 593 to be detected and avoided.
- 594 • If the duration distributions are not exponential, we only get unbiased  
 595 predictors, by using the proposed approach with truncated distribu-  
 596 tions.

597 Note that we did not compare to prediction approaches clustering cases  
 598 by system load [9], or using queueing models [35]. This is intentional, as  
 599 these further conditions are orthogonal to the time domain and could be  
 600 integrated with the approach presented in this work.

### 601 5.6. Scalability Analysis

602 We propose to use simulation to predict the remaining service execution  
 603 time for a running case. Therefore, it is interesting to see how long the  
 604 prediction approach takes to simulate the remaining behavior of the process.  
 605 To test for scalability, we conduct the following experiment.

606 First, we randomly generate structured, acyclic GDT\_SPN models of dif-  
 607 ferent sizes by iterative insertion of sequential, parallel, and exclusive blocks,  
 608 until we reach the desired node size of the net. For each timed transition

609 we specify a distribution, i.e., either *uniform*, *normal*, or *lognormal* distri-  
610 butions, or non-parametric *Gaussian kernel* density estimators based on a  
611 random number of observations. In general, the run-time of the prediction  
612 algorithm depends on multiple factors:

- 613 • The *average number of transitions* that need to fire to reach the end  
614 of the process, influenced by the size of the net, the progress of the  
615 current case, and potential cycles.
- 616 • The *kind of transition distributions*, as there exist distributions that are  
617 rather costly to draw samples from, e.g., complicated non-parametric  
618 models, as well as very simple models, e.g., the uniform distribution.
- 619 • The *requested accuracy of the prediction*. Besides the fact that comput-  
620 ing a narrow confidence interval takes more samples than allowing more  
621 sampling error, the variance of the process durations also influences the  
622 number of samples required to achieve the requested precision.
- 623 • The *computing power* of the system running the simulation.

624 For our experiments, we fixed the *requested accuracy* to a 99 percent confi-  
625 dence interval within maximum  $\pm 3$  percent of error of the mean remaining  
626 duration. We control for the variance of the activity durations when assigning  
627 the distributions. Regarding *computing power*, we used a desktop computer  
628 with a Pentium i5-4570 (3.20 GHz cores) equipped with 16GB of ram. We  
629 varied the other two factors, i.e., the *average number of transitions* and the  
630 *kinds of distributions* used. Figure 8 depicts the average time taken for re-  
631 maining time prediction of acyclic GDT.SPN models based on the number of  
632 transitions in the model in log-scale. For example, predicting the duration  
633 of a model with 300 transitions takes well below 100 milliseconds for all dis-  
634 tributions. Note that a prediction of a model with 10000 transitions still is  
635 feasible in less than 100 seconds with these configurations. We can see that  
636 the computation time scales almost linearly with the size of the model.

637 In our experience, most business processes involving human activities take  
638 hours, days, or sometimes even months to complete. In these situations, the  
639 quality of the prediction is more important than the performance of the  
640 prediction approach. If performance is critical however, the approach could  
641 be extended to provide a fall back option to an analytical method based on  
642 GSPN models, as provided by Zimmermann [40].



## 643 **6. Conclusion**

644 This paper presented an approach for predicting the remaining time, and  
645 for estimating the chance of missing a given deadline. The approach is appli-  
646 cable with a limited number of assumptions and yields significant improve-  
647 ments over related work by conditioning on elapsed time. The effects of  
648 conditioning on time are most prominent when there are long time spans  
649 between observable events in a business process. The underlying model is  
650 able to reflect concurrency in processes better than the state of the art [8]  
651 and its derivatives, as it builds on the Petri net formalism.

652 Future work includes taking resource dependencies, or dependencies be-  
653 tween activity durations into account. In many cases, process data plays  
654 an important role for the remaining durations of activities, and often affects  
655 decisions in the process. Adding data dependencies to the model, as for ex-  
656 ample investigated in [38], is complementary to the proposed approach and  
657 promising as well.

## 658 **References**

- 659 [1] M. Weske, *Business Process Management: Concepts, Languages, Archi-*  
660 *tectures*, 2nd Edition, Springer, 2012.
- 661 [2] M. Dumas, M. La Rosa, J. Mendling, H. A. Reijers, *Fundamentals of*  
662 *Business Process Management*, Springer, 2013.
- 663 [3] B. Mutschler, M. Reichert, J. Bumiller, Unleashing the effectiveness of  
664 process-oriented information systems: Problem analysis, critical success  
665 factors, and implications, *Systems, Man, and Cybernetics, Part C: Ap-*  
666 *plications and Reviews*, IEEE Transactions on 38 (3) (2008) 280–291.
- 667 [4] S. Taylor, Waiting for service: The relationship between delays and  
668 evaluations of service, *Journal of Marketing* 58 (2) (1994) 56–69.
- 669 [5] P. Kumar, M. U. Kalwani, M. Dada, The impact of waiting time guaran-  
670 tees on customers’ waiting experiences, *Marketing Science* 16 (4) (1997)  
671 295–314.
- 672 [6] S. Meyer, A. Ruppen, C. Magerkurth, Internet of things-aware process  
673 modeling: Integrating iot devices as business process resources, in: *Ad-*  
674 *vanced Information Systems Engineering*, Vol. 7908 of LNCS, Springer,  
675 2013, pp. 84–98.

- 676 [7] A. Rogge-Solti, Probabilistic estimation of unobserved process events,  
677 Ph.D. thesis, Hasso Plattner Institute at the University of Potsdam,  
678 Germany (April 2014).
- 679 [8] W. M. P. van der Aalst, M. H. Schonenberg, M. Song, Time prediction  
680 based on process mining, *Information Systems* 36 (2) (2011) 450–475.
- 681 [9] F. Folino, M. Guarascio, L. Pontieri, Discovering Context-Aware Models  
682 for Predicting Business Process Performances, in: *OTM 2012*, Springer,  
683 2012, pp. 287–304.
- 684 [10] A. Rogge-Solti, M. Weske, Prediction of remaining service execution  
685 time using stochastic petri nets with arbitrary firing delays, in: *ICSOC*,  
686 Vol. 8274 of LNCS, Springer, Heidelberg, 2013, pp. 389–403.
- 687 [11] R. Conforti, M. de Leoni, M. La Rosa, W. M. P. van der Aalst, Support-  
688 ing risk-informed decisions during business process execution, in: *Ad-  
689 vanced Information Systems Engineering*, Vol. 7908 of LNCS, Springer,  
690 2013, pp. 116–132.
- 691 [12] A. Pika, W. M. P. van der Aalst, C. J. Fidge, A. H. M. ter Hofstede,  
692 M. T. Wynn, Predicting deadline transgressions using event logs, in:  
693 *BPM Workshops 2012*, Vol. 132 of LNBIP, Springer, 2013, pp. 211–216.
- 694 [13] B. Denton, J. Viapiano, A. Vogl, Optimization of surgery sequencing  
695 and scheduling decisions under uncertainty, *Health Care Management  
696 Science* 10 (1) (2007) 13–24.
- 697 [14] O. M. G. (OMG), Business Process Model and Notation (BPMN) 2.0  
698 Specification (January 2011).
- 699 [15] C. A. Petri, Kommunikation mit Automaten, Ph.D. thesis, Technische  
700 Hochschule Darmstadt (1962).
- 701 [16] W. M. P. van der Aalst, The application of petri nets to workflow man-  
702 agement, *Journal of Circuits, Systems, and Computers* 8 (1) (1998) 21–  
703 66.
- 704 [17] N. Lohmann, H. M. W. E. Verbeek, R. Dijkman, Petri Net Transfor-  
705 mations for Business Processes - A Survey, in: *Transactions on Petri  
706 Nets and Other Models of Concurrency II*, Vol. 5460 of LNCS, Springer,  
707 2009, pp. 46–63.

- 708 [18] W. M. P. van der Aalst, Verification of workflow nets, in: ICATPN'97,  
709 Vol. 1248 of LNCS, Springer, 1997, pp. 407–426.
- 710 [19] F. J. W. Symons, Modeling and Analysis of Communication Protocols  
711 Using Numerical Petri Nets, Ph.D. thesis, University of Essex, Great  
712 Britain (1978).
- 713 [20] M. A. Marsan, G. Conte, G. Balbo, A Class of Generalized Stochastic  
714 Petri Nets for the Performance Evaluation of Multiprocessor Systems,  
715 ACM TOCS 2 (2) (1984) 93–122.
- 716 [21] S. Natkin, Reseaux de Petri Stochastiques, Ph.D. thesis, Conservatoire  
717 National des Arts et Métiers (CNAM), Paris, France (1980).
- 718 [22] M. K. Molloy, Performance Analysis Using Stochastic Petri Nets, Com-  
719 puters, IEEE Transactions on 100 (9) (1982) 913–917.
- 720 [23] A. Horváth, A. Puliafito, M. Scarpa, M. Telek, Analysis and Evaluation  
721 of Non-Markovian Stochastic Petri Nets, in: Computer Performance  
722 Evaluation. Modelling Techniques and Tools, Springer, 2000, pp. 171–  
723 187.
- 724 [24] A. Rogge-Solti, W. M. P. van der Aalst, M. Weske, Discovering stochas-  
725 tic petri nets with arbitrary delay distributions from event logs, in: BPM  
726 Workshops, Vol. 171 of LNBIP, Springer, Heidelberg, 2014, pp. 15–27.
- 727 [25] B. F. van Dongen, A. K. A. de Medeiros, H. M. W. Verbeek, A. J.  
728 Weijters, W. M. P. van der Aalst, The ProM framework: A new era in  
729 process mining tool support, in: Applications and Theory of Petri Nets  
730 2005, Springer, 2005, pp. 1105–1116.
- 731 [26] B. F. van Dongen, R. A. Crooy, W. M. P. van der Aalst, Cycle time  
732 prediction: When will this case finally be finished?, in: OTM 2008,  
733 Springer, 2008, pp. 319–336.
- 734 [27] W. M. P. van der Aalst, A. Adriansyah, B. F. van Dongen, Replaying  
735 history on process models for conformance checking and performance  
736 analysis, in: WIREs: Data Mining and Knowledge Discovery, Vol. 2,  
737 Wiley Online Library, 2012, pp. 182–192.

- 738 [28] W. J. Dixon, F. J. Massey, et al., Introduction to statistical analysis,  
739 Vol. 344, McGraw-Hill New York, 1969.
- 740 [29] S. F. Chen, J. Goodman, An Empirical Study of Smoothing Techniques  
741 for Language Modeling, in: Proceedings of the 34th annual meeting on  
742 Association for Computational Linguistics, Association for Computa-  
743 tional Linguistics, 1996, pp. 310–318.
- 744 [30] S.-Y. Hwang, H. Wang, J. Tang, J. Srivastava, A Probabilistic Approach  
745 to Modeling and Estimating the QoS of Web-Services-Based Workflows,  
746 Information Sciences 177 (23) (2007) 5484–5503.
- 747 [31] H. Zheng, J. Yang, W. Zhao, A. Bouguettaya, QoS Analysis for Web  
748 Service Compositions Based on Probabilistic QoS, in: Service-Oriented  
749 Computing 2011, Springer, 2011, pp. 47–61.
- 750 [32] B. Kiepuszewski, A. H. M. t. Hofstede, C. J. Bussler, On Structured  
751 Workflow Modelling, in: Advanced Information Systems Engineering,  
752 Vol. 1789 of LNCS, Springer, 2000, pp. 431–445.
- 753 [33] A. Wombacher, M.-E. Iacob, Estimating the Processing Time of Pro-  
754 cess Instances in Semi-structured Processes—A Case Study, in: Services  
755 Computing (SCC), 2012 IEEE Ninth International Conference on, IEEE,  
756 2012, pp. 368–375.
- 757 [34] K. P. Anklesaria, Z. Drezner, A multivariate approach to estimating the  
758 completion time for PERT networks, The Journal of the Operational  
759 Research Society 37 (8) (1986) 811–815.
- 760 [35] A. Senderovich, M. Weidlich, A. Gal, A. Mandelbaum, Queue mining—  
761 predicting delays in service processes, in: CAiSE, Springer, 2014, pp.  
762 42–57.
- 763 [36] M. de Leoni, W. M. van der Aalst, M. Dees, A general framework for cor-  
764 relating business process characteristics, in: BPM, Vol. 8659 of LNCS,  
765 Springer, 2014, pp. 250–266.
- 766 [37] A. Rozinat, M. T. Wynn, W. M. P. van der Aalst, A. H. M. ter Hofstede,  
767 C. J. Fidge, Workflow simulation for operational decision support, Data  
768 & Knowledge Engineering 68 (9) (2009) 834–850.

- 769 [38] A. Rozinat, R. S. Mans, M. Song, W. M. P. van der Aalst, Discovering  
770 simulation models, *Information Systems* 34 (3) (2009) 305–327.
- 771 [39] G. Balbo, G. Chiola, Stochastic Petri Net Simulation, in: *Proceedings*  
772 *of the 21st conference on Winter simulation*, ACM, 1989, pp. 266–276.
- 773 [40] A. Zimmermann, Modeling and Evaluation of Stochastic Petri Nets with  
774 TimeNET 4.1, in: *VALUETOOLS’12*, IEEE, 2012, pp. 54–63.
- 775 [41] M. Rosemann, M. zur Mühlen, Integrating risks in business process mod-  
776 els, in: *ACIS 2005 Proceedings*, 2005, paper 50.
- 777 [42] S. Suriadi, B. Weiss, A. Winkelmann, A. H. ter Hofstede, M. Adams,  
778 R. Conforti, C. Fidge, M. L. Rosa, C. Ouyang, M. Rosemann, A. Pika,  
779 M. Wynn, Current research in risk-aware business process management:  
780 overview, comparison, and gap analysis, *Communications of the Asso-*  
781 *ciation for Information Systems* 34 (1) (2014) 933–984.
- 782 [43] G. Goluch, A. Ekelhart, S. Fenz, S. Jakoubi, S. Tjoa, T. Mück, Integra-  
783 tion of an ontological information security concept in risk aware business  
784 process management, in: *HICSS-41*, 2008, pp. 377–386.
- 785 [44] B. Kang, D. Kim, S.-H. Kang, Periodic Performance Prediction for Real-  
786 time Business Process Monitoring, *Industrial Management & Data Sys-*  
787 *tems* 112 (1) (2012) 4–23.
- 788 [45] P. Leitner, B. Wetzstein, F. Rosenberg, A. Michlmayr, S. Dustdar,  
789 F. Leymann, Runtime Prediction of Service Level Agreement Viola-  
790 tions for Composite Services, in: *ICSOC/ServiceWave 2009 Workshops*,  
791 Springer, 2010, pp. 176–186.
- 792 [46] W. M. P. van der Aalst, T. Weijters, L. Maruster, Workflow mining:  
793 Discovering process models from event logs, *Knowledge and Data Engi-*  
794 *neering*, *IEEE Transactions on* 16 (9) (2004) 1128–1142.
- 795 [47] S. J. J. Leemans, D. Fahland, W. M. van der Aalst, Discovering block-  
796 structured process models from event logs containing infrequent be-  
797 haviour, in: *BPM Workshops*, Vol. 171 of LNBIP, Springer, 2014, pp.  
798 66–78.
- 799 [48] B. F. v. Dongen, BPI Challenge 2012 Dataset, [http://dx.doi.org/10.](http://dx.doi.org/10.4121/uuid:3926db30-f712-4394-aebc-75976070e91f)  
800 [4121/uuid:3926db30-f712-4394-aebc-75976070e91f](http://dx.doi.org/10.4121/uuid:3926db30-f712-4394-aebc-75976070e91f) (2012).