

Exploring LLM Features in Predictive Process Monitoring for Small-Scale Event-Logs

ALESSANDRO PADELLA, Università degli Studi di Padova, Italy

MASSIMILIANO DE LEONI, Università degli Studi di Padova, Italy

MARLON DUMAS, University of Tartu, Estonia

Predictive Process Monitoring is a branch of process mining that aims to predict the outcome of an ongoing process. Recently, it leveraged machine-and-deep learning architectures. In this paper, we extend our prior LLM-based Predictive Process Monitoring framework, which was initially focused on total time prediction via prompting. The extension consists of comprehensively evaluating its generality, semantic leverage, and reasoning mechanisms, also across multiple Key Performance Indicators. Empirical evaluations conducted on three distinct event logs and across the Key Performance Indicators of Total Time and Activity Occurrence prediction indicate that, in data-scarce settings with only 100 traces, the LLM surpasses the benchmark methods. Furthermore, the experiments also show that the LLM exploits both its embodied prior knowledge and the internal correlations among training traces. Finally, we examine the reasoning strategies employed by the model, demonstrating that the LLM does not merely replicate existing predictive methods but performs higher-order reasoning to generate the predictions.

Additional Key Words and Phrases: Predictive process monitoring, Large language models, Trace Encoding

ACM Reference Format:

Alessandro Padella, Massimiliano de Leoni, and Marlon Dumas. 2026. Exploring LLM Features in Predictive Process Monitoring for Small-Scale Event-Logs. 1, 1 (January 2026), 19 pages.

1 Introduction

Predictive Process Monitoring (PPM) is a family of techniques that leverages event logs from business processes to generate predictions about the future states or properties of ongoing process instances [17]. PPM methods vary depending on the prediction target, which can include times [34], next activities [24], or process outcomes [32].

Literature has extensively explored machine-and-deep learning models to enhance prediction quality [6]. However, these models typically require large amounts of data for effective training. When the available event log is limited in size, the applicability of such techniques becomes constrained, reducing the overall potential of PPM. As highlighted in [37], data availability remains one of the most significant challenges faced by researchers and practitioners in this domain. Large Language Models (LLMs) provided an alternative for the application of PPM to data-scarce environments since their embedded knowledge enables robust prediction, while structured prompting preserves trace sequences, dependencies, and attributes without extensive fine-tuning [15, 29].

This paper extends our prior LLM-based PPM work [25]. Our previous work focused on total time prediction via Gemini¹ prompting, which showed to be successful in this context when the training consisted of event logs of limited size, it remained unclear whether the prediction quality would extend to other Key Performance Indicators (KPIs)

¹<https://docs.cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-flash?hl=en>

Authors' Contact Information: [Alessandro Padella](mailto:alessandro.padella@unipd.it), alessandro.padella@unipd.it, Università degli Studi di Padova, Padova, Italy; [Massimiliano de Leoni](mailto:deleoni@math.unipd.it), deleoni@math.unipd.it; [Marlon Dumas](mailto:marlon.dumas@ut.ee), University of Tartu, Tartu, Estonia, marlon.dumas@ut.ee.

whose predictions are relevant in different contexts. Our previous work also fell short in a comprehensive analysis of the mechanisms and knowledge used behind to enhance the prediction quality. In this paper, we aim to extend the prior work investigating three research questions.

- (1) **RQ1:** When trained on event logs with a limited number of traces, do LLMs achieve superior prediction quality on a wide range of metrics, compared to the models available in literature?
- (2) **RQ2:** Do LLMs leverage embodied prior knowledge of the process domains to improve predictions?
- (3) **RQ3:** Do LLMs use an ensemble of different prediction models? If so, how would this take place?

To address RQ1, we extend the experimental scenario to predicting whether certain activities will occur. All results are statistically validated and benchmarked against state-of-the-art models, including CatBoost[9] and PGTNet[10].

To address RQ2, we introduce semantic hashing to enhance the prompt formulation. Specifically, we encode all process-related strings, such as trace variable and activity names, into hashed representations. We then repeat the experiments using these modified prompts. This procedure enables us to assess whether the LLM’s performance relies solely on event correlations or also leverages contextual information.

To answer RQ3, we distill the reasoning processes of the LLM into interpretable β -learners derived through pattern analysis of 150 traces. These learners are subsequently re-implemented and evaluated to determine whether they can match or exceed the performance of the original LLM. Additionally, we conduct statistical significance tests to ensure the robustness and consistency of the derived β -learners.

The remainder of this paper is organized as follows: Section 2 reviews literature related to LLM in the field of process mining. Section 3 presents the preliminaries needed to define the approach in Section 4, in which the prompting technique is reported and detailed. Finally, Section 5 reports the extended version of experiments, while Section 6 concludes the paper.

2 Related Works

LLMs are recently attracting growing focus in business process management as well [2]. They have proven to be significantly useful for many process mining tasks, such as process modeling [3, 18], log extraction [7], anomaly detection [33], and they have also been used for assessing the validity of some new given traces for a given process [28].

Lashkevich et al. in [21] provide a state-of-the-art approach that leverages LLMs for enhancing the optimization of waiting times and relies on user-prompted feedback for recommending more effective re-design options. Rebmman et al. [28] present an approach for extracting knowledge from textual data, providing textual and synthetically generated benchmark datasets for extracting event logs to assess missing activities and generate them. The work in [26] leverages LLMs to transform textual data into process representations, followed by training a BERT-based deep learning model to predict the next activity in a process.

Berti et al. [3] fine-tuned a pretrained LLM using reinforcement learning to generate complete and executable process models from textual descriptions, employing structural and behavioral rewards to improve correctness and reduce invalid generations. Casciani et al. [5] introduce a retrieval-augmented LLM framework for next activity prediction in predictive process monitoring using past event traces without training, while identifying limits such as interleaving sensitivity and concept drift. The approach in [14] contextually links event logs with additional process-related data from varied sources, enabling LLMs to provide relevant process insights via natural language querying for planning, monitoring, improving operations and anomaly detection [16]. Finally, Kubrak et al. [19] developed a chatbot-based approach for process analysis where the LLM explains recommendations from a model to enhance explainability.

3 Preliminaries

The starting point for a process mining-based system is an *event log*. An event log is a multiset of *traces*. Each trace is a sequence of events, each describing a particular *process instance* (i.e., a *case*) in terms of the *activities* executed, the associated *timestamps* and other different domain-related *attributes*.

Definition 3.1 (Events). Let \mathcal{A} be the set of process activities. Let \mathcal{T} be the set of process timestamps. Let $\mathcal{V} = \mathcal{V}_1 \times \mathcal{V}_2 \times \dots \times \mathcal{V}_m$ be the Cartesian product of the data attribute sets. An event is a tuple $(a, t_{start}, t_{end}, \vec{v}) \in \mathcal{A} \times \mathcal{T}^2 \times \mathcal{V}$ where a is the event activity, t_{start} and t_{end} the associated timestamps, and \vec{v} the vector of associated attributes.

A trace is a sequence of events. The same event can occur in different traces. Namely, attributes may be given the same assignment in different traces. This means that the same trace can appear multiple times, although admittedly under extremely rare conditions, and motivates why an event log has to be defined as a multiset of traces:

Definition 3.2 (Traces & Event Logs). Let $\mathcal{E} = \mathcal{A} \times \mathcal{T}^2 \times \mathcal{V}$ be the universe of events. A trace σ is a sequence of events, i.e. $\sigma \in \mathcal{E}^*$.² An event log \mathcal{L} is a multiset of traces, i.e. $\mathcal{L} \subset \mathbb{B}(\mathcal{E}^*)$.³

Given an event $e = (a, t_{start}, t_{end}, \vec{v})$, the remainder uses the following shortcuts: $activity(e) = a$, $start(e) = t_{start}$, $end(e) = t_{end}$, $duration(e) = t_{start} - t_{end}$, $attr(e) = \vec{v}$. Also, given a single attribute set \mathcal{V}_i , it is associated an attribute name, i.e. $name(\mathcal{V}_i)$ and it can be classified as *global* or *local*, i.e. $type(\mathcal{V}_i) \in \{global, local\}$ depending on whether the values in it can vary or not in the same trace. We refer to the value of these attributes as $global(\sigma) = \vec{g}$ and $local(e) = \vec{l}$, and so the equation $global(\sigma) \oplus local(e) = attr(e)$ holds.⁴ Furthermore, given a trace $\sigma = \langle e_1, \dots, e_n \rangle$, $prefix(\sigma)$ denotes the set of all prefixes of σ , including σ , namely $prefix(\sigma) = \{\langle \rangle, \langle e_1 \rangle, \langle e_1, e_2 \rangle, \dots, \langle e_1, \dots, e_n \rangle\}$.

The goal of a KPI prediction framework is to forecast the KPI value of a process instance that has not completed yet, namely a *running trace*.

In this paper, the problem is modeled as the estimation of a KPI function $\mathcal{K} : \mathcal{X} \rightarrow \mathbb{R}_0$ that given a running trace $\sigma' = \langle e_1, \dots, e_k \rangle$ eventually completing as $\langle e_{k+1}, \dots, e_n \rangle$, returns the KPI value $k \in \mathbb{R}_0$ after the occurrence of all the events $\langle e_1, \dots, e_n \rangle$ in the trace. The input of the KPI function is a set \mathcal{X} , since not every approach shares the same encoding for event logs. For instance, in [1] the authors encoded traces in an LSTM compatible input, while in [30] the traces are encoded in a *comma-separated values* file suitable for a predictor based on a Decision Tree. This requires defining the **trace-to-instance encoding function** $\rho : \mathcal{E}^* \rightarrow \mathcal{X}$ with the goal of accurately translating every trace of the event log into an input suitable for the predictive model. This function has proven to be significantly different based on the chosen predictive approach (cf. [31]).

Figure 1 depicts an example of a trace-to-instance encoding function. In it, referred as ρ_{aggr} the trace is preprocessed as a row of a .csv file by adding the past activities in newly generated columns. For each activity in the trace, the number of previous occurrences of that activity is reported in a dedicated column, encoding the number of past executions of the activity. This encoding allows tracking the frequency of all past activities but does not maintain information about their sequential order, recording only the most recent one.

²The operator $*$ refers to the Kleene star: given a set A , A^* contains all the possible finite sequences of elements belonging to A .

³ $\mathbb{B}(X)$ indicates the set of all multisets with the elements in set X .

⁴Considering \oplus as the concatenation of vectors e.g.

$[1, 3, 'request_created'] \oplus [2, True] = [1, 3, 'request_created', 2, True]$

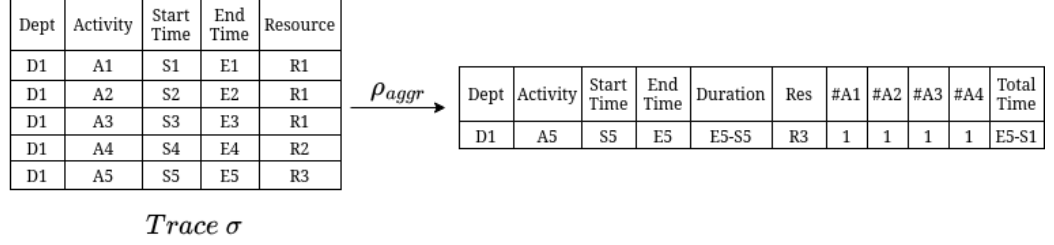


Fig. 1. Example of an encoding function applied to a running trace used for predicting the KPI “Total Time”.

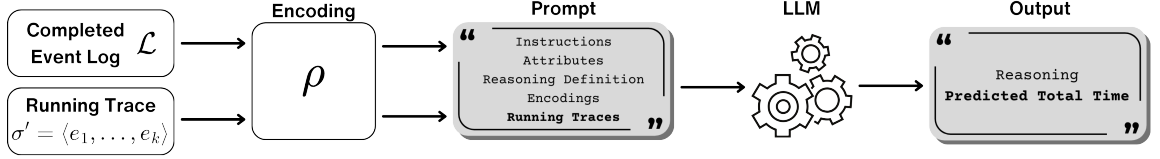


Fig. 2. Pipeline outlining the proposed method using LLMs for PPM, in which the KPI is “Total Time”.

4 Approach For LLM-based Predictions

This study seeks to leverage the potential of LLMs to develop a framework for PPM, particularly in scenarios where only a small amount of example traces are accessible. Leveraging their embedded knowledge, LLMs can extract and use additional information beyond the event data by incorporating the semantics of events, such as activity names, that traditional models cannot. Figure 2 depicts the proposed approach. Given an event log of completed traces \mathcal{L} and a running trace σ' , a trace-to-instance encoding function ρ is applied to transform them into a structured prompt. This prompt, composed of multiple components, is then used to enable the LLM to estimate the KPI function \mathcal{K} .

In the remainder of this Section a new trace-to-instance encoding function ρ_{seq} suitable for LLM is introduced in Section 4.1, while Section 4.2 defines a context-based prompt suitable for employing an LLM for implementing the KPI function defined in Section 3, while Section 4.3 provides an example of LLM’s output containing the predicted values and its reasoning to achieve it.

4.1 An Encoding Function for LLMs

Exploiting an LLM for developing a PPM framework is a topic that has not yet been explored in process mining (cf. Section 2). This section introduces a new LLM-suitable trace-to-instance encoding function $\rho_{seq} : \mathcal{E}^* \rightarrow \Sigma^*$.⁵

The input of the ρ_{seq} function is a trace, while the output is a textual prompt that will be later enhanced to become suitable as prompt for an LLM, that will be used as KPI function \mathcal{K} . Specifically, in this case the generic input set of the KPI function \mathcal{X} is equal to Σ^* . In ρ_{seq} , each trace $\sigma = \langle e_1, \dots, e_n \rangle$ is mapped into a string composed of three main elements:

- The values of the global attributes of the trace *global* (e_n).
- A sequence of tuples (*activity* (e_i), *duration* (e_i)) for $i = 1, \dots, n$.
- The actual value of $\mathcal{K}(\sigma)$.

⁵Considering Σ as the set of all finite strings over an alphabet.

```
" { Dept : D1, ActTimeSeq : { (A1, E1-S1), (A2, E2-S2), (A3, E3-S3), (A4, E4-S4), (A5, E5-S5) }, TotalTime : E5-S1} "
```

Fig. 3. Example of a usage of the Sequential encoding function ρ_{aggr} applied to the same running trace and KPI employed in Figure 1.

Formally:

$$\rho_{seq}(\langle e_1, \dots, e_n \rangle) = global(\langle e_1, \dots, e_n \rangle) \oplus (activity(e_1), duration(e_1)) \oplus \dots \oplus (activity(e_n), duration(e_n)) \oplus (k)$$

The sets of local attributes have been intentionally excluded, as it has been demonstrated that LLMs are constrained by two primary factors: **technical** limitations and **methodological** considerations. From a technical perspective, an LLM can only process a certain number of characters; so it becomes necessary to reduce the size of the input to stay within this maximum quantity, namely the **Context Length**.⁶ Note that the Context Length of an LLM is not just a limitation per interaction (e.g., in a chatbot) but an inherent architectural constraint. Additionally, from a methodological standpoint, research demonstrates that the data in an LLM input does not equally impact the model's processing, and the significance of individual data points reduces as the input lengthens, even degrading its performance [20, 22]. Therefore, we opted to omit local attributes. Conversely, global attributes were retained since they incorporate domain knowledge and have proved to retain more predictive power than local ones (cf. Galanti et al. [12]).

An example of application of the Sequential trace-to-instance encoding function ρ_{seq} is depicted in Figure 3, where the example trace reported in Figure 1 has been processed as a string. The result is a string-form Python object primarily composed of three keys: i) *Dept*, associated with the value of the corresponding global attribute; ii) *ActTimeSeq*, associated with the list of tuples where each activity and its duration are recorded; and iii) *TotalTime*, representing the total duration of the trace.

It is important to note that the Sequential is not the only suitable encoding for LLMs. In fact, the proposed approach also supports the *Aggregated History Encoding*, depicted in Figure 1, as the LLM can predict future values independently of the specific encoding. This flexibility allows the framework to adapt to different encoding strategies, that have been explored in [25].

4.2 Context-Based Prompting Technique for LLMs in Predictive Process Monitoring

This Section uses the traces that have been encoded using a trace-to-instance encoding function ρ and incorporates them generating an input suitable for an LLM. In essence we define a prompting technique that allows the model to generate KPI predictions along with corresponding reasoning procedure, starting from the encoded traces. The prompting technique is divided into seven key parts, also reported in the Listing 1:

- **Initial instruction and Header:** The LLM is introduced to the task with the prompt: “You are an expert in process mining and machine learning. Your task is to predict the KPI of process instances based on event logs of activities, where each process instance is a sequence of activities.” (Lines 1-2)
- **Attributes and Encoding description:** Contextual information specific to the process is provided, and the trace-to-instance encoding function is described (ρ_{seq} in the example). (Lines 4-12)
- **Output and Reasoning format specification:** The expected structure for predicted values is defined. (Lines 14-22 and 29-35)

⁶See <https://llm-stats.com/> for an overview of Context Lengths of the latest models.

- **Running Trace Format Specification:** The format for describing a running trace is specified to the model. (Lines 24-27)
- **Domain-specific background information:** Additional details about the process from which the data have been extracted. (Lines 37-38)
- **Example Data Provision:** Encoded data are provided as example to the model. (Lines 39-45)
- **Running Trace Provision:** The running trace is provided in the same format as the examples, with a custom last activity referred as “Running”, as described at lines 24-27. (Lines 47-50)

Although the proposed encoding is general and applicable to various use cases and KPI, certain information within these seven components must be specified by the process analyst and may be optionally removed.

```

1 You are an expert in process mining and machine learning. Your task is to predict the 'total time' of
2 process instances based on event logs, as each process instance is a sequence of activities.
3
4 A event log is a collection of traces, where each trace represents a process instance.
5 Each trace is mapped as a sequence of activities and integers representing the minutes since the start
6 of the process.
7 The log is represented as a python list containing one dictionary for each trace. Included in it are:
8 - the key "AMOUNT_REQ", representing the total amount of euros requested in the loan application.
9 - the key "ActTimeSeq", which value is a list of [activity, cumulative elapsed minutes]
10 - The key "total_time", which value is the total execution time in minutes from the start of the activity,
11 that is the value to predict.
12
13 All interactions will be structured in the following way, with the appropriate values filled in.
14
15 [[ ## reasoning ## ]]
16 {your step-by-step reasoning}
17
18 [[ ## answer ## ]]
19 {your predicted total time as an integer}
20
21 [[ ## completed ## ]]
22
23 In adhering to this structure, your objective is to analyze the event log, and apply reasoning to predict
24 the total time for a new case. This case belongs to a not-yet-completed process instance, represented by the
25 label "Running" in "ActTimeSeq", indicating that more activities are expected before reaching the conclusion
26 of the process instance.
27
28 Ensure to articulate each step of your thought process in the reasoning field, detailing how you identify
29 relationships with past cases and leverage your intuition about the meaning of activities to arrive at the
30 solution. The answer should be the final prediction of the total time for the given process instance.
31 Respond with the corresponding output fields, starting with the field [[ ## reasoning ## ]],
32 then [[ ## answer ## ]], and then ending with the marker for [[ ## completed ## ]].
33
34 Your task is to learn from them and predict the 'total time' values for that traces.
35
36 The process deals with a loan application process from a Dutch financial institution. It has been provided
37 in the Business Process Intelligence (BPI) challenge in 2012.
38 The following list shows some completed example cases with their total times:

```

```

39
40     {"AMOUNT_REQ": 5000.0, "ActTimeSeq": [["W_Completeren aanvraag", 11],
41     ["W_Nabellen offertes", 1464], ["W_Nabellen offertes", 7486]], "total_time": "7486"}
42     {"AMOUNT_REQ": 15000.0, "ActTimeSeq": [["W_Completeren aanvraag", 13],
43     ["W_Nabellen offertes", 14], ["W_Validate application", 4328], ["W_Validate application", 8792]],
44     "total_time": "8792"}
45
46 Now predict the total time for this new uncompleted case, considering that the case is still running:
47
48     {"Application_1000386745": {"AMOUNT_REQ": 18000.0,
49     "ActTimeSeq": [["W_Completeren aanvraag", 2], ["W_Nabellen offertes", 8571], ["Running"]]}

```

Listing 1. Prompting technique example for a loan application process (Bpi12), for the KPI of Total Time. Lines that have to be provided by the process analyst are marked in bold. In the example, only 2 training traces are provided due to space limitation.

In Listing 1 they have been highlighted in bold. They are specifically, (i) the domain-specific background information (Lines 37-38) and (ii) the description of the global attributes (Lines 8-9), as they contain contextual details specific to the process under study. Notably, these informations are optional and can be excluded if necessary. The remaining sections of the prompt are designed to be generic and can be applied to any event log without modification.

This modularity ensures that the framework can be associated with any encoding function and eventually be customized with details about the single process, while minimizing the effort needed for customization when applying the LLM to different process datasets.

4.3 Example of Reasoning to Achieve Predictions

The LLM's output is not limited to returning the predicted value: it also reports on the reasoning put in place to conclude the prediction. This is in fact a form of explanation of the predictions, of which one example is given in Listing 2.

```

1  [[ ## reasoning ## ]]
2  The new case is for an application with a 'AMOUNT_REQ' of 18000. The process instance
3  has executed activities 'W_Completeren aanvraag' and 'W_Nabellen offertes'.
4  The process is currently 'Running'.
5  To predict the total time, we will look at historical cases with similar characteristics.
6
7  1. Filter for similar Application Type and Requested Amount:
8      We focus on cases with 'AMOUNT_REQ'
9      in a range around 18000.0 (e.g., 10000.0 to 25000.0) to find comparable instances.
10
11  2. Identify relevant historical cases:
12      From the provided log, we select the cases that match 'AMOUNT_REQ' within the range of 10000.0 to 25000.0.
13      We list their total times:
14      10000.0: 9713, 10000.0: 18803, 11500.0: 4771, 12000.0: 18924, 12000.0: 8102,
15      13000.0: 6645, 13000.0: 21550, 14000.0: 9713, 15000.0: 8049, 15000.0: 17011
16
17  3. Calculate the median total time:
18      To get a robust estimate and reduce the impact of outliers, we calculate the median of these total times.
19      Sorted total times: 4771, 5299, 5975, 5975, 6645, 6645, 7432, 8049, 8102, 9713, 9713,
20      12991, 13056, 17011, 17434, 18803, 18924, 19826, 21550, 71762.
21      71762 is detected as an outlier and it will be removed from the analysis.

```

```

22     Since there are 19 values, the median is the 10th value,
23     which is 9713. Median = 9713.
24
25 4. Predict total time:
26     We use the median total time of the similar historical cases as our prediction
27     for the new running case.
28
29 [[ ## answer ## ]]
30 9713
31 [[ ## completed ## ]]

```

Listing 2. LLM’s output and reasoning returned after the prompting technique example for a loan application process.

Nonetheless, it is a valuable addition to this paper to report and comment on an instance of the prediction’s reasoning given as output by the LLM. To do so, Example 2 reports the output returned by the LLM for Example 1.

The output is returned at the end (see line 32 of Example 2). The reasoning of this example is as follows: the LLM has focused on the traces of the provided training log in which the request amount is in the range of 10000 to 25000 (see lines 7-9). This leads to 20 traces (see lines 15-16), for which the median value is computed (see lines 24-25), which is used as predicted value (see explanations given at lines 28 and 29).

This reasoning certainly provides a valid explanation of the reasoning behind why this prediction is provided. However, this is only an example of the reasoning procedure performed by the LLM, and we observed that the output prompt can vary on the basis of the context and the number of training traces provided.

5 Experimental Results

We assess the effectiveness of our approach in generating valuable predictions using a limited-size training set and compare its performance against state-of-the-art benchmarks across three distinct use cases on two different KPIs. To ensure a robust evaluation, we repeatedly sampled a limited number of traces, trained both a benchmark model and the LLMs on these samples, and conducted multiple experimental runs. The remainder of this section is organized as follows: Section 5.1 reports how the logs have been divided into training and test sets for use in the experimental setup. The use cases are presented in Section 5.2, while Section 5.3 outlines the experimental pipeline employed to address the research questions posed in Section 1. Finally, Section 5.4 concludes with the results and their associated analysis.

5.1 LLM used and Train and Test Split

The whole approach has been implemented in Python and the code is publicly available.⁷ The prediction function development has been carried on using **Gemini 2.5 Flash Thinking**: a state-of-the-art LLM developed by Google DeepMind⁸, but any choice of LLM is valid. The model is built on a multimodal architecture designed for advanced natural language understanding and generation. As the development of LLMs progresses, we anticipate that the performances of upcoming models using our method may see fast improvements.

Consistently with standard supervised learning practices, we divided the event log \mathcal{L} into training and test, \mathcal{L}^{comp} and \mathcal{L}^{run} , respectively. To extract the training log we compute the earliest time t_{split} such that 80% of the identifiers related to traces of \mathcal{L} are completed. This allows us to define \mathcal{L}^{comp} as the set of traces of \mathcal{L} completed at time t_{split} , and consequently, define \mathcal{L}^{run} as $\mathcal{L} \setminus \mathcal{L}^{comp}$. The traces of the test log \mathcal{L}^{run} are truncated to a set \mathcal{L}^{trunc} , namely the

⁷https://github.com/Pado123/gui_xrecs_presc_analytics

⁸<https://gemini.google.com/app>

| Use Case | Training Traces | 100/Training_Traces | Target Activity |
|----------|-----------------|---------------------|--|
| Bpi12 | 6792 | 1.45% | W_Nabellen incomplete dossiers |
| Bac | 25843 | 0.38% | Service closure Request with BO responsibility |
| Hospital | 30394 | 0.32% | LABORATORIO |

Table 1. Overview of the analyzed use cases, including the number of available training traces, the relative proportion represented by 100 traces, and the target activity selected for the KPI Activity Occurrence.

set of prefixes, that is obtained from \mathcal{L}^{run} by removing every event with a timestamp larger than t_{split} : \mathcal{L}^{trunc} only contains the events that occurred before time t_{split} . This procedure tries to mimic the reality at time t_{split} and it is in line with the principles introduced in [35]. The system is trained on \mathcal{L}_{comp} , the predictions are produced for \mathcal{L}^{trunc} and tested using its completed form, \mathcal{L}_{run} . Furthermore, to ensure a robust generalization across various process instances and a more balanced comparison with the LLM, we randomly picked 10% of traces and used them as a validation set \mathcal{L}_{valid} to apply a Cross-Validation approach to optimize the following parameters of the benchmarks. This technique is a widespread technique in PPM works [4, 8, 27]

It is also worthwhile pointing out that in this paper we use the term *training* to refer to both LLMs and machine- and deep-learning benchmark models, to keep the discussion simple. However, we acknowledge that LLMs are already pre-trained: traces are provided to the LLMs as background, and are not formally used to train its internal parameters.

5.2 Use Cases and Benchmarks

The evaluation considers three real-life event logs that are commonly used in PPM:

- **Bpi12** This process has been used by the BPI challenge in 2012⁹, it contains 8,616 traces, 6 different activities and 1 global attribute: *Requested_Amount*.
- **Bac** A process referring to a process of a Bank Institution that deals with the closures of bank accounts. It contains 32,429 completed traces, 15 different activities and 2 global attributes: *Closure_Type*, and *Closure_Reason*.¹⁰
- **Hospital** This process has been provided by an hospital emergency department. The log is made of 37,945 completed traces, contains 46 different activities and 3 global attributes, that are *Triage_Color*, *Triage_Access* and *Patient_Age*. Due to confidentiality, this event log cannot be published.

For each use case, two prediction tasks, corresponding to two different KPIs:

- **Total Time** of an ongoing case, i.e., a regression task.
- **Activity Occurrence** for a selected target class, i.e. a classification task. The chosen

Table 1 summarizes the analyzed use cases, reporting the proportion that 100 traces represent relative to the total number of training traces. The last column identifies the target activity for the KPI prediction task of Activity Occurrence. These target activities were selected because their execution typically leads to higher time and cost demands. In the Bpi12 use case, the target activity is *W_Nabellen incomplete dossiers*, which translates to “Follow up on incomplete files.” This activity occurs when an error in the loan application process requires rework by both the sides of the bank and the customer. Similarly, the activities “Service Closure with BO Responsibility” and “Laboratorio”, that indicates the requirements of further blood analysis, implying additional laboratory processing and resource consumption.

⁹<https://doi.org/10.4121/uuid:3926db30-f712-4394-aebc-75976070e91f>

¹⁰https://github.com/IBM/processmining/tree/main/Datasets_usecases

For the purpose of setting a benchmark for prediction, we employed two different State-of-the-art techniques for the different KPIs of interest:

- For Total Time we employed PGTNet architecture from [10], that has proven to outperform every framework.
- For the Activity Occurrence, we employed Catboost from [9], a state-of-the-art model predictor based on machine learning on decision trees, which has been shown to surpass existing prediction frameworks [11].

5.3 Experimental Pipeline

With the extent of answering to the three research questions posed in Section 1, the experimental pipeline comprises three different main steps:

5.3.1 Prediction Quality Comparison. This first step of the pipeline extends the experimental part of [25] and addresses RQ1: *When trained on event logs with a limited number of traces, do LLMs achieve superior prediction quality on a wide range of metrics, compared to the models available in literature?*

We measured the accuracy of the LLM when trained on 100 running traces sampled from \mathcal{L}^{train} , that are a limited amount of training traces with respect to the model being trained on the whole training set (cf. Table 1). To evaluate the performance of the prediction models, two key metrics are employed depending on the KPI of interest:

- **Mean Absolute Error (MAE)** For Total Time prediction, the Mean Absolute Error is used:

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

where y_i is the actual Total Time, reported in minutes, \hat{y}_i is the predicted value, and n is the number of samples. Lower MAE values indicate better predictive performance.

- **F1-Score** For activity occurrence prediction, the F1-Score is used:

$$F1\text{-Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

where Precision = $\frac{TP}{TP+FP}$ and Recall = $\frac{TP}{TP+FN}$, with TP , FP , and FN denoting true positives, false positives, and false negatives.

Since this 100 traces have been randomly sampled for each run, we repeated the experiments 20 times in order to asses statistical validity, this motivates why the resulting values of MAE/F1-Score are reported associated with the corresponding standard deviation.

5.3.2 Hashing of the Input Prompts. The second step of the experimental pipeline addresses RQ2: *Do LLMs leverage embodied prior knowledge of the process domains to improve predictions?*

To isolate semantic understanding from distributional patterns, we anonymize all context-sensitive strings that could enable the LLM to exploit domain knowledge. Traditional machine-and-deep learning methods remain unaffected by string anonymization since they operate on numerical encodings, conversely LLMs with embodied knowledge may leverage activity semantics. For instance, the activity name “LABORATORIO” implies laboratory analysis, but also and attribute names as, “Triage_Color” suggests emergency care. This leads to the fact that we have to anonymize not only the values of the categorical attributes and activities, but also the name of that attributes since they provide information about the process’ context.

We define the context-sensitive set \mathcal{H} as:

$$\mathcal{H} = \mathcal{A} \cup \left(\bigcup_{i=1}^n \{v \in \mathcal{V}_i \mid \text{type}(\mathcal{V}_i) = \text{global}\} \right) \cup \left(\bigcup_{i=1}^n \{\text{name}(\mathcal{V}_i) \mid \text{type}(\mathcal{V}_i) = \text{global}\} \right) \quad (1)$$

A deterministic hash function $H : \mathcal{H} \rightarrow \Sigma^4$ maps each $s \in \mathcal{H}$ to a unique 4-character identifier ($c_1 c_2 c_3 c_4$, $c_i \in \Sigma = \text{A-Z,0-9}$), preserving correlations while eliminating semantics.

In this experimental pipeline’s phase, after the prompt generation as outlined in Section 4, every $s \in \mathcal{H}$ is replaced by $H(s)$. Then, the prediction quality in terms of MAE or F1-Score is compared between hashed (anonymized) and non-hashed prompts to quantify the semantic dependency of the LLM. Significant degradation under hashing indicates that the LLM is relying on its embodied knowledge.

Furthermore, we aim to assess the statistical significance between hashed and non-hashed prompts. To this extent we applied the Nemenyi post-hoc test across repeated experiments:

The Nemenyi test, from Nemenyi [23], is a global statistical test (such as the Friedman test) that serves to investigate if it is possible to reject the null hypothesis that the performance of the comparisons on the groups of data is similar. The test makes pair-wise tests of performance and we apply it under the following hypothesis:

- Null hypothesis (H_0): No significant difference exists between prediction performance (MAE/F1-Score) of hashed and non-hashed prompts across use cases.
- Alternative hypothesis (H_1): Significant performance difference exists, with non-hashed prompts expected to outperform due to semantic exploitation.

5.3.3 Derivation and Re-implementation of LLM’s β -learners. The third step of our experimental pipeline addresses RQ3: *Do LLMs use an ensemble of different prediction models? If so, how would this take place?*

To investigate this question, we manually analyzed 50 LLM reasoning explanations per use case and KPI, yielding a total of 150 reasoning traces per KPI. Following the visualization methodology proposed in [36], we systematically cataloged recurring decision patterns extracted from these textual explanations. Specifically, the aim of this part is to identify how the LLM justifies its predictions as requested in the prompt reported in the Listing 1.

These recurring patterns were abstracted into families, the so-called β -learners. In this context a β -learners constitutes a mathematically defined and reproducible model that explicitly performs the reasoning strategies that the LLM documented .

For each use case, every derived β -learner is applied to the same test set \mathcal{L}^{test} introduced in 5.1. Then, we compared prediction quality in terms of MAE/F1-Score against the LLM results.

Furthermore, the statistical significance of performance differences is assessed via Wilcoxon signed-rank test ($\alpha = 0.05$) with the following hypotheses for each β -learner vs. LLM comparison with the following hypothesis:

- **Null hypothesis (H_0):** No significant difference exists between prediction performance (MAE/F1-Score) of the current β -learner and the one obtained by the LLM.
- **Alternative hypothesis (H_1):** A significant difference exists between prediction performance (MAE/F1-Score) of the current β -learner and the one obtained by the LLM.

Furthermore, to determine whether the LLM’s predictions merely replicate β -learner strategies or exhibit more sophisticated inference, we compared each β -learner’s prediction quality **as employed within the LLM’s own reasoning** against the LLM’s overall results. In a nutshell, we report the percentage difference in terms of MAE/F1-Score between the LLM and each β -learner when the LLM has been classified using that β -learner.

Table 2. MAE (minutes) for hashed/non-hashed prompts

| Use Case | Model | hash | all_df | 100 ± |
|----------|--------------------|------|--------|-------------|
| Bpi12 | CatBoost Regressor | - | 6846 | 9394 ± 1275 |
| | PGTNet | - | 3888 | 8856 ± 576 |
| | 2.5 flash | no | - | 6508 ± 235 |
| | 2.5 flash | yes | - | 9246 ± 873 |
| Bac | CatBoost Regressor | - | 2647 | 6393 ± 365 |
| | PGTNet | - | 1245 | 4753 ± 472 |
| | 2.5 flash | no | - | 2265 ± 1072 |
| | 2.5 flash | yes | - | 3880 ± 3254 |
| hospital | CatBoost Regressor | - | 253 | 259 ± 52 |
| | PGTNet | - | 97 | 132 ± 32 |
| | 2.5 flash | no | - | 115 ± 34 |
| | 2.5 flash | yes | - | 2077 ± 232 |

Finally, we want to understand if 150 LLM reasoning per KPI were enough to derive all the β -learners. To this extent, we employ *Good-Turing frequency estimation*. This statistical technique, originally developed for natural language processing and machine translation, provides a principled framework for estimating the probability of a new unseen event based on the frequency distribution of observed ones [13].

The Good-Turing formula estimates the adjusted probability of an item that has been observed r times as:

$$P^*(r) = \frac{(r+1) \cdot N_{r+1}}{N_r \cdot N} \quad (2)$$

where:

- N_r denotes the number of β -learner patterns that have been observed exactly r times in the training set of 50 derived patterns,
- N_{r+1} denotes the number of patterns observed exactly $(r+1)$ times,
- N is the total number of pattern occurrences,
- $P^*(r)$ is the smoothed probability estimate.

The probability of observing a *new* (never before seen) β -learner is estimated as:

$$P_0 = \frac{N_1}{N} \quad (3)$$

where N_1 is the number of β -learner patterns observed exactly once in the training set.

This approach allows us to compute the expected number of novel β -learners $E[\text{novel}]$ that would be encountered when applying the method to m new prediction instances (where $m \in \{1, 10, 100\}$):

$$E[\text{novel}] = m \cdot P_0 = m \cdot \frac{N_1}{N} \quad (4)$$

5.4 Results

5.4.1 Quality Prediction and Hashing Results. Table 2 presents MAE in minutes for the prediction of the KPI Total Time across Bpi12, Bac, and Hospital use cases, comparing with the benchmark of CatBoost Regressor and PGTNet, as

Table 3. F1-Score for hashed/non-hashed prompts

| Use Case | Model | hash | all_df | 100 ± |
|----------|---------------------|------|--------|-------------|
| Bpi12 | CatBoost Classifier | - | 0.80 | 0.72 ± 0.07 |
| | 2.5 flash | no | - | 0.77 ± 0.06 |
| | 2.5 flash | yes | - | 0.74 ± 0.06 |
| Bac | CatBoost Classifier | - | 0.95 | 0.78 ± 0.08 |
| | 2.5 flash | no | - | 0.98 ± 0.04 |
| | 2.5 flash | yes | - | 0.96 ± 0.04 |
| hospital | CatBoost Classifier | - | 0.90 | 0.90 ± 0.01 |
| | 2.5 flash | no | - | 0.90 ± 0.08 |
| | 2.5 flash | yes | - | 0.84 ± 0.06 |

Table 4. Summary of Nemenyi Post-Hoc Test Results

| Task | Use Case | Difference | P-value | Significance |
|---------------------|----------|------------|---------|--------------|
| Total Time | Bpi12 | 42.06% | 0.002 | ** |
| | Bac | 71.25% | <0.001 | *** |
| | Hospital | 1702% | 0.002 | ** |
| Activity Occurrence | Bpi12 | 0.03 | 0.018 | * |
| | Bac | 0.02 | 0.040 | * |
| | Hospital | 0.06 | <0.001 | *** |

introduced in Section 1. Results also reports the MAE when the benchmarks used all training data (all_df), scarce-data (100 traces, $\tilde{1}$ -0.3% of logs per Table 1), and hashed vs. non-hashed, associated with the standard deviation.

When trained on 100 traces, non-hashed LLM yields competitive MAE: Bpi12 (6508 ± 235) surpasses CatBoost (9394 ± 1275) and nears PGTNet (8856 ± 576) while Hospital (115 ± 34) outperforms both (259 ± 52 , 132 ± 32). In a similar manner Bac (2265 ± 1072) outperforms both PGTNet (4753 ± 472) and CatBoost.

Hashing activity/attribute semantics (e.g., masking "LABORATORIO") degrades LLM's performance of Bpi12 to 9246 ± 873 (+42%), Hospital to 2077 ± 232 (+1700%), Bac to 3880 ± 3254 .

Table 3 reports F1-scores for predicting targeted activity occurrence, comparing CatBoost Classifier as benchmark. As for Table 2, the results for the benchmark have been also reported when it has been trained on the whole \mathcal{L}^{train}

When trained on 100 traces, non-hashed LLM achieves comparable or superior performance with respect to the benchmark: Bpi12 (0.77 ± 0.06 vs. 0.72 ± 0.07), Bac (0.98 ± 0.04 markedly exceeding 0.78 ± 0.08), Hospital (0.90 ± 0.08 matching 0.90 ± 0.01). Furthermore, low standard deviations confirm prediction stability, demonstrating LLMs' capacity to forecast activity occurrences despite minimal training traces. Addressing RQ1.

Hashing yields consistent yet moderate degradation, such as Bpi12 to 0.74 ± 0.06 (-4%), Bac to 0.96 ± 0.04 (-2%), Hospital to 0.84 ± 0.06 (-7%).

Table 4 outlines the results on the Nemenyi tests for addressing RQ2. All the comparisons reject the null hypothesis of no difference between hashed and non-hashed prompts, confirming the LLM's reliance on semantic embodied knowledge. For total time prediction, effect sizes are substantial, also rising to 1702% difference in Hospital use case. These findings, aligned with MAE/F1 trends in Tables 2 and 3, demonstrate that LLMs exploit event semantics.

5.4.2 β -Learner Performance Analysis. For the prediction of the KPI Total Time, the extracted β -learners are:

- **knn act + mean/median/mode**: k -nearest neighbors on activity-based representations with different aggregation functions.
- **knn att + mean/median/mode**: k -nearest neighbors on attribute-based representations with analogous aggregations functions.
- **time seq + mean/median/mode**: models that exploit the temporal sequence of events with different aggregation functions.
- **path pred + mean/median/mode**: models based on predicted future paths of a case, again with different aggregations functions.

For prediction of the KPI activity occurrence prediction, different β -learners have been identified:

- **Activity-Based**: Based on analyzing the sequence of activities
- **State-Based**: Based on analyzing the last event
- **Att-Based**: Based on analyzing the trace attributes values.
- **Positive Evidence**: Checking if the activity to prediction already occurred or not.

The details on how the single β -learners have been implemented are reported in the code url.

Tables 5 and 3 present the MAE and F1-Score results describing the β -learners' predictive performance compared to the LLM. The tables systematically compare the re-implementations of the LLM's reasoning patterns, i.e., the β -learners such as knn-act/att-mean/mode (nearest neighbors on activity/attribute features), time-seq-mean (temporal sequence aggregation), and path-pred-mean/mode (predicted future path estimation), across all three use cases under both the full \mathcal{L}^{train} and when trained on 100 traces.

Table 5 reports almost all significantly different results (cf. Significance column), showing that the LLM consistently outperforms every associated β -learner with strong statistical support. This trend is slightly different in Table 3, where, for the Bpi12 use case, the p -value indicating the difference is less significant than in the others. However, the results remain stable and statistically meaningful.

In the Δ LLM column, the percentage difference in MAE and F1-Score relative to the LLM is reported for Table 5 and 3, respectively. Since every result is far from 0, we can conclude that the LLM model does not simply replicate the β -learners, but instead performs substantially better.

5.4.3 Good-Turing Results. Table 7 provides results in terms of *the probability of encounter a new β -learner* analyzing other 1, 10 or 100 LLM's outputs. Strong empirical evidence supporting the claim that the probability of discovering novel β -learner patterns in future predictions approaches zero. Expected novel β -learners for $m = 1, 10, 100$ additional traces approach zero in every use case and KPI. Specifically, after 1 new running trace, the probability is exactly 0 for both KPI of Total Time and Activity Occurrence. After 10 new running traces, it raises to 0.1% and 0.4% for respectively Total Time and Activity Occurrence, reaching 0.3% and 1.4% after 100 new running traces.

Furthermore, Figure 4 plots the convergence curves of the different KPI distributions (Total Time on top rows and Activity Occurrence on bottom rows) when analyzing a given number of traces, also divided per use case. From the different plot it is possible to visualize that the convergence in the distribution has been already reached after roughly 30 traces.

| Use Case | Model | all_df | 100 ± | #best50 no | Significance | Occurrence | ΔLLM when occurring |
|----------|-------------------|--------|--------------|------------|--------------|------------|---------------------|
| Bpi12 | 2.5 flash no-hash | - | 6508 ± 235 | 39 | - | - | - |
| | knn act mean | 8356 | 8651 ± 393 | 2 | *** | 5 | 6% |
| | knn act median | 8082 | 8780 ± 524 | 2 | *** | 5 | 8% |
| | knn act mode | 13463 | 13200 ± 1002 | 0 | *** | 6 | 9% |
| | knn att mean | 13261 | 8974 ± 245 | 1 | *** | 4 | 75% |
| | knn att median | 14253 | 8957 ± 453 | 0 | *** | 2 | 51% |
| | knn att mode | 14591 | 16077 ± 1690 | 0 | *** | 4 | 24% |
| | time seq mean | 8764 | 9127 ± 537 | 1 | *** | 3 | 26% |
| | time seq median | 8016 | 8614 ± 772 | 4 | *** | 6 | 7% |
| | time seq mode | 12433 | 12710 ± 1636 | 0 | *** | 5 | 9% |
| | path pred mean | 15023 | 15060 ± 3087 | 0 | *** | 4 | 23% |
| | path pred median | 12349 | 13468 ± 2819 | 1 | *** | 4 | 23% |
| | path pred mode | 8546 | 15998 ± 4578 | 0 | *** | 6 | 11% |
| Bac | 2.5 no-hash | - | 2265 ± 874 | 36 | - | - | - |
| | knn act mean | 3994 | 6788 ± 3972 | 1 | *** | 4 | 24% |
| | knn act median | 3563 | 3095 ± 437 | 4 | ** | 9 | 10% |
| | knn act mode | 3372 | 5503 ± 3434 | 0 | ** | 6 | 22% |
| | knn att mean | 21296 | 8094 ± 2874 | 0 | *** | 3 | 41% |
| | knn att median | 5401 | 4678 ± 1969 | 1 | ** | 3 | 22% |
| | knn att mode | 5379 | 14964 ± 879 | 1 | *** | 3 | 11% |
| | time seq mean | 6917 | 9224 ± 3820 | 0 | *** | 2 | 8% |
| | time seq median | 6966 | 3916 ± 531 | 1 | *** | 5 | 12% |
| | time seq mode | 8702 | 9534 ± 4135 | 0 | *** | 1 | 22% |
| | path pred mean | 54692 | 6843 ± 3820 | 0 | *** | 6 | 61% |
| | path pred median | 3569 | 3249 ± 445 | 4 | *** | 9 | 8% |
| | path pred mode | 190092 | 5927 ± 2699 | 0 | *** | 4 | 6% |
| Hospital | 2.5 flash no-hash | - | 115 ± 60 | 45 | - | - | - |
| | knn act mean | 167 | 442 ± 379 | 0 | ** | 4 | 24% |
| | knn act median | 388 | 175 ± 15 | 1 | ** | 7 | 11% |
| | knn act mode | 458 | 503 ± 653 | 0 | *** | 2 | 12% |
| | knn att mean | 617 | 566 ± 437 | 0 | *** | 3 | 12% |
| | knn att median | 189 | 183 ± 27 | 1 | *** | 5 | 80% |
| | knn att mode | 1479 | 626 ± 546 | 0 | *** | 7 | 68% |
| | time seq mean | 661 | 670 ± 805 | 0 | *** | 4 | 11% |
| | time seq median | 684 | 287 ± 677 | 0 | *** | 2 | 61% |
| | time seq mode | 643 | 708 ± 771 | 1 | *** | 4 | 34% |
| | path pred mean | 383 | 382 ± 320 | 0 | *** | 4 | 32% |
| | path pred median | 159 | 174 ± 15 | 2 | ** | 4 | 14% |
| | path pred mode | 282 | 473 ± 618 | 0 | *** | 7 | 31% |

Table 5. Complete MAE results across all use cases and models, also paired with the possible outcome of β -learners.

6 Conclusion

This paper advances our previous work [25] in the field of Predictive Process Monitoring in data-scarce environments through a comprehensive LLM prompting framework. It rigorously addresses three new research questions via empirical analysis across three public and private event logs and two different Key Performance Indicators (KPIs).

RQ1 confirms the LLM’s superiority when the model is trained on 100 traces, outperforming state-of-the-art benchmarks for both KPIs. This demonstrates that the LLM is a viable alternative for implementing Predictive Process Monitoring frameworks.

RQ2 empirically proves semantic leverage: the context in the prompts was hashed and the experiments repeated, showing that the LLM exploits its embodied knowledge to provide predictions for running traces. The results were also statistically verified using the Nemenyi post-hoc test.

RQ3 reveals the reasoning anatomy: manual inspection of 50 traces for each use case and KPI led to the identification of different predictors that the LLM employs, the so-called β -learners. These have been shown to include all possible

Table 6. F1-Score results related to β -learners.

| Use Case | Model | hash | all_df | 100 \pm | Significance | Occurrence | Δ LLM when occurring |
|----------|-------------------|------|--------|-----------------|--------------|------------|-----------------------------|
| Bpi12 | 2.5 flash | no | - | 0.77 ± 0.06 | - | - | - |
| | Activity-Based | - | 0.71 | 0.71 ± 0.05 | ** | 15 | 0.03 |
| | State-Based | - | 0.70 | 0.70 ± 0.02 | * | 14 | 0.02 |
| | Positive Evidence | - | 0.57 | 0.57 ± 0.01 | ** | 17 | 0.07 |
| | Att-Based | - | 0.50 | 0.50 ± 0.10 | ** | 12 | 0.06 |
| Bac | 2.5 flash | no | - | 0.98 ± 0.04 | - | - | - |
| | Activity-Based | - | 0.78 | 0.75 ± 0.04 | *** | 10 | 0.04 |
| | State-Based | - | 0.72 | 0.72 ± 0.05 | *** | 21 | 0.24 |
| | Positive Evidence | - | 0.58 | 0.58 ± 0.04 | *** | 10 | 0.25 |
| | Att-Based | - | 0.52 | 0.52 ± 0.03 | *** | 9 | 0.19 |
| Hospital | 2.5 flash | no | - | 0.90 ± 0.08 | - | - | - |
| | Activity-Based | - | 0.82 | 0.80 ± 0.05 | ** | 15 | 0.09 |
| | State-Based | - | 0.75 | 0.71 ± 0.04 | *** | 10 | 0.18 |
| | Positive Evidence | - | 0.58 | 0.58 ± 0.04 | *** | 12 | 0.14 |
| | Att-Based | - | 0.53 | 0.53 ± 0.06 | *** | 14 | 0.25 |

Table 7. Good-Turing Smoothing Analysis: Novel β -Learner Discovery Rate

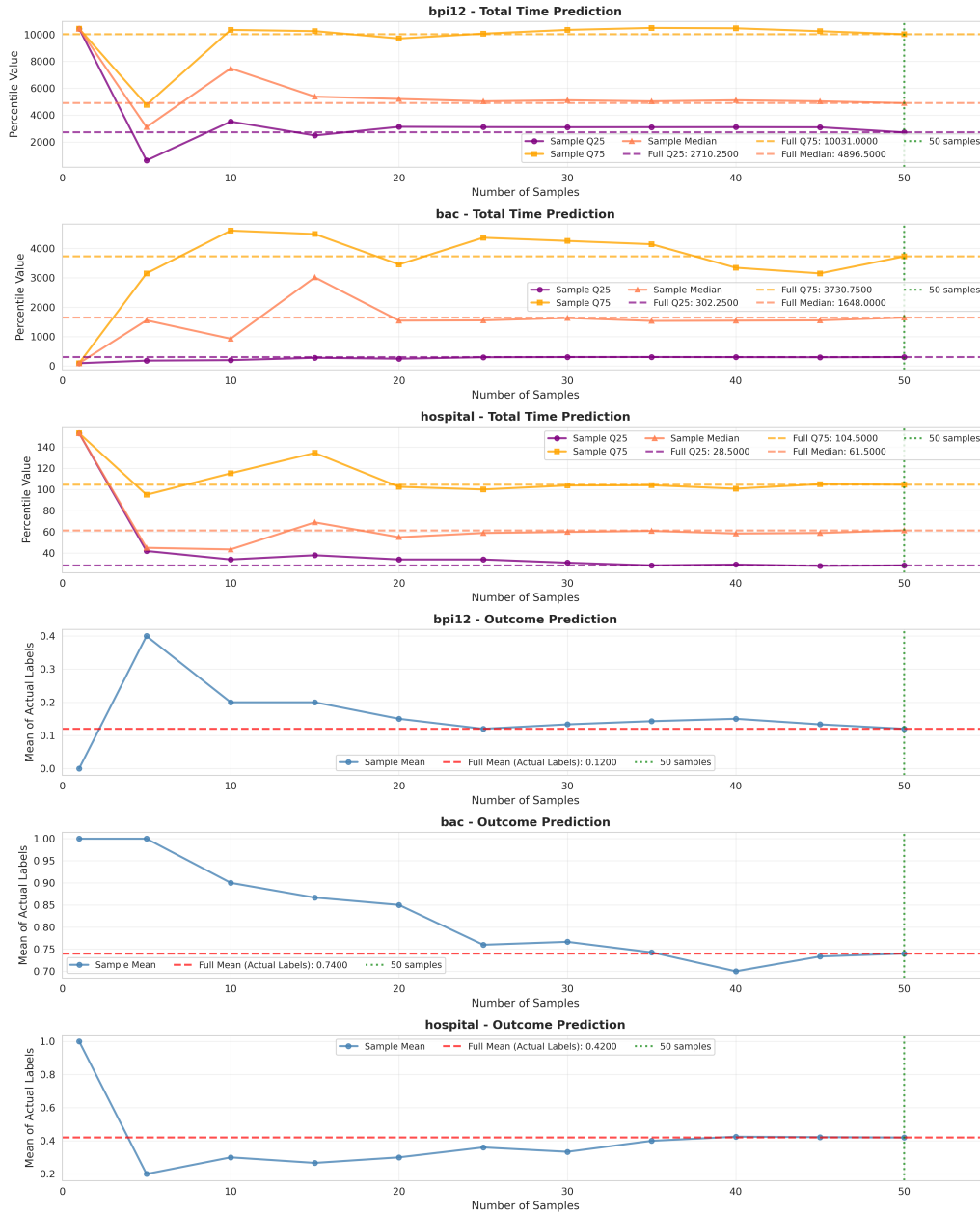
| Task | Expected Novel β -Learners |
|---------------------|---|
| Total Time | $m = 1: 0.000, m = 10: 0.001, m = 100: 0.003$ |
| Activity Occurrence | $m = 1: 0.000, m = 10: 0.004, m = 100: 0.014$ |

ones through Good-Turing validation. This paper demonstrates not only that the LLM outperforms every model it mimics but also that it performs more complex analyses than those models.

Future work will focus on exploring new LLMs, given the rapid evolution of the field, and will include user studies on the reliance on β -learners. Another direction is to extend this work toward a Prescriptive Process Analytics framework, enabling not only predictions but also actionable recommendations derived from them.

Fig. 4. KPI convergences of various use cases.

Convergence Analysis: All Case Studies (Non-Hashed, n100)
Total Time: Q25, Q75, Median | Outcome: Mean



References

- [1] Muhammad Awais Ali, Marlon Dumas, and Fredrik Milani. 2024. Enhancing the Accuracy of Predictors of Activity Sequences of Business Processes. In *Research Challenges in Information Science*. Springer Nature Switzerland.

- [2] Alessandro Berti, Humam Kourani, and Wil M. P. van der Aalst. 2025. PM-LLM-Benchmark: Evaluating Large Language Models on Process Mining Tasks. In *Process Mining Workshops*, Andrea Delgado and Tijs Slaats (Eds.). Springer Nature Switzerland, Cham, 610–623.
- [3] Alessandro Berti, Xiaoting Wang, Humam Kourani, and Wil M. P. Van der Aalst. 2025. Specializing large language models for process modeling via reinforcement learning with verifiable and universal rewards. *Process Science* 2, 1 (2025), 26. doi:10.1007/s44311-025-00034-4
- [4] Andrei Buliga, Chiara Di Francescomarino, Chiara Ghidini, Ivan Donadello, and Fabrizio Maria Maggi. 2025. Guiding the generation of counterfactual explanations through temporal background knowledge for predictive process monitoring. *Data Mining and Knowledge Discovery* 39, 5 (2025), 63. doi:10.1007/s10618-025-01117-3
- [5] Angelo Casciani, Mario Luca Bernardi, Marta Cimitile, and Andrea Marrella. 2026. Enhancing next activity prediction in process mining with Retrieval-Augmented Generation. *Information Systems* 137 (2026), 102642. doi:10.1016/j.is.2025.102642
- [6] Paolo Ceravolo, Marco Comuzzi, Jochen De Weerdt, et al. 2024. Predictive Process Monitoring: Concepts, Challenges, and Future Research Directions. *Process Science* 1, 2 (2024), 2.
- [7] Vinicius Stein Dani, Marcus Dees, Henrik Leopold, Kiran Busch, Iris Beerepoot, Jan Martijn E. M. van der Werf, and Hajo A. Reijers. 2024. Event Log Extraction for Process Mining Using Large Language Models. In *Cooperative Information Systems - 30th International Conference, CoopIS 2024, Porto, Portugal (Lecture Notes in Computer Science)*. Springer.
- [8] Massimiliano de Leoni and Alessandro Padella. 2024. Achieving Fairness in Predictive Process Analytics via Adversarial Learning. In *Cooperative Information Systems - 30th International Conference, CoopIS 2024, Porto, Portugal, November 19-21, 2024, Proceedings (Lecture Notes in Computer Science, Vol. 15506)*, Marco Comuzzi, Daniela Grigori, Mohamed Sellami, and Zhangbing Zhou (Eds.). Springer, 346–354. doi:10.1007/978-3-031-81375-7_21
- [9] Anna Veronika Dorogush, Vasily Ershov, and Andrey Gulin. 2017. CatBoost: gradient boosting with categorical features support. In *Proceedings of the Workshop on ML Systems at NIPS 2017*.
- [10] Keyvan Amiri Elyasi, Han van der Aa, and Heiner Stuckenschmidt. 2024. PGTNet: A Process Graph Transformer Network for Remaining Time Prediction of Business Process Instances. In *Advanced Information Systems Engineering - 36th International Conference, CAISE 2024, Limassol, Cyprus, 2024, Proceedings (Lecture Notes in Computer Science, Vol. 14663)*. Springer, 124–140.
- [11] Riccardo Galanti, Bernat Coma-Puig, Massimiliano de Leoni, Josep Carmona, and Nicolò Navarin. 2020. Explainable Predictive Process Monitoring. In *Proceedings of the 2nd International Conference on Process Mining (ICPM 2020)*. IEEE.
- [12] Riccardo Galanti, Massimiliano de Leoni, Merylin Monaro, Nicolò Navarin, Alan Marazzi, Brigida Di Stasi, and Stéphanie Maldera. 2023. An explainable decision support system for predictive process analytics. *Eng. Appl. Artif. Intell.* 120 (2023), 105904.
- [13] I. J. Good. 1953. The population frequencies of species and the estimation of population parameters. *Biometrika* 40, 3-4 (1953), 237–264.
- [14] Sathyanarayanan Gopalan, Vijay V, Sreenidhi Kalyanaraman, Daniel Joseph Raj, Sneha V, and Anto N Thomas. 2024. Leveraging LLMs for Context-Aware Process Mining and Analysis on Multi-Source Data. In *2024 IEEE 8th International Conference on Information and Communication Technology (CICT)*. 1–6. doi:10.1109/CICT64037.2024.10899487
- [15] Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew G Wilson. 2023. Large language models are zero-shot time series forecasters. *Advances in Neural Information Processing Systems* 36 (2023), 19622–19635.
- [16] Sascha Kaltenpoth, Alexander Skolik, Oliver Müller, and Daniel Beverungen. 2026. A Step Towards Cognitive Automation: Integrating LLM Agents with Process Rules. In *Business Process Management, Arik Senderovich, Cristina Cabanillas, Irene Vanderfeesten, and Hajo A. Reijers (Eds.)*. Springer Nature Switzerland, Cham, 308–324.
- [17] Jongchan Kim, Marco Comuzzi, Marlon Dumas, Fabrizio Maria Maggi, and Irene Teinemaa. 2022. Encoding resource experience for predictive process monitoring. *Decision Support Systems* 153 (2022).
- [18] Humam Kourani, Alessandro Berti, Daniel Schuster, and Wil M. P. van der Aalst. 2024. Process Modeling with Large Language Models. In *Enterprise, Business-Process and Information Systems Modeling*. Springer Nature Switzerland, Cham, 229–244.
- [19] Kateryna Kubrak, Lana Botchorishvili, Fredrik Milani, Alexander Nolte, and Marlon Dumas. 2024. Explanatory Capabilities of Large Language Models in Prescriptive Process Monitoring. In *Business Process Management - 22nd International Conference, BPM 2024, Krakow , 2024, Proceedings (Lecture Notes in Computer Science, Vol. 14940)*. Springer, 403–420.
- [20] Yury Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. 2024. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack. *Advances in Neural Information Processing Systems* 37 (2024), 106519–106554.
- [21] Katsiaryna Lashkevich, Fredrik Milani, Maksym Avramenko, and Marlon Dumas. 2024. LLM-Assisted Optimization of Waiting Time in Business Processes: A Prompting Method. In *Business Process Management - 22nd International Conference, BPM 2024, Krakow , 2024, Proceedings (Lecture Notes in Computer Science, Vol. 14940)*. Springer, 474–492.
- [22] Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhui Chen. 2024. Long-context LLMs Struggle with Long In-context Learning. *Transactions on Machine Learning Research* (2024).
- [23] Peter B. Nemenyi. 1963. *Distribution-free Multiple Comparisons*. Ph. D. Dissertation. Princeton University.
- [24] Rafael Oyamada, Gabriel Tavares, Sylvio Barbon Junior, and Paolo Ceravolo. 2024. *Enhancing Predictive Process Monitoring with Time-Related Feature Engineering*. Springer Nature Switzerland.
- [25] Alessandro Padella, Paolo Frazzetto, Nicolò Navarin, and Massimiliano de Leoni. 2026. Enhancing Predictive Process Monitoring on Small-Scale Event Logs Using LLMs. In *Business Process Management Forum, Arik Senderovich, Cristina Cabanillas, Irene Vanderfeesten, and Hajo A. Reijers (Eds.)*. Springer Nature Switzerland, Cham, 274–290.

- [26] Vincenzo Pasquadibisceglie, Annalisa Appice, and Donato Malerba. 2024. LUPIN: A LLM Approach for Activity Suffix Prediction in Business Process Event Logs. In *2024 6th International Conference on Process Mining (ICPM)*. 1–8.
- [27] Jari Peepkorn, Seppe vanden Broucke, and Jochen De Weerd. 2024. Validation set sampling strategies for predictive process monitoring. *Information Systems* 121 (2024), 102330. doi:10.1016/j.is.2023.102330
- [28] Adrian Rebmann, Fabian David Schmidt, Goran Glavas, and Han van der Aa. 2024. Evaluating the Ability of LLMs to Solve Semantics-Aware Process Mining Tasks. In *6th International Conference on Process Mining, ICPM 2024, Kgs. Lyngby, Denmark, October 14-18, 2024*. IEEE, 9–16.
- [29] James Requeima, John Bronskill, Dami Choi, Richard E. Turner, and David Kristjanson Duvenaud. 2024. LLM Processes: Numerical Predictive Distributions Conditioned on Natural Language. In *Advances in Neural Information Processing Systems*.
- [30] Mahmoud Shoush and Marlon Dumas. 2025. White box specification of intervention policies for prescriptive process monitoring. *Data Knowl. Eng.* 155 (2025), 102379.
- [31] Gabriel M. Tavares, Rafael Seidi Oyamada, Sylvio Barbon, and Paolo Ceravolo. 2023. Trace encoding in process mining: A survey and benchmarking. *Eng. Appl. Artif. Intell.* 126, Part D (2023), 107028.
- [32] Irene Teinemaa, Marlon Dumas, Marcello La Rosa, and Fabrizio Maggi. 2017. Outcome-Oriented Predictive Process Monitoring: Review and Benchmark. *ACM Transactions on Knowledge Discovery from Data* 13 (07 2017).
- [33] Han van der Aa, Adrian Rebmann, and Henrik Leopold. 2021. Natural language-based detection of semantic execution anomalies in event logs. *Information Systems* 102 (2021), 101824.
- [34] Ilya Verenich, Marlon Dumas, Marcello La Rosa, Fabrizio Maggi, and Irene Teinemaa. 2019. Survey and Cross-benchmark Comparison of Remaining Time Prediction Methods in Business Process Monitoring. *ACM Transactions on Intelligent Systems and Technology* 10 (07 2019).
- [35] Hans Weytjens and Jochen De Weerd. 2021. Creating Unbiased Public Benchmark Datasets with Data Leakage Prevention for Predictive Process Monitoring. In *Business Process Management Workshops - BPM 2021 International Workshops, Rome, Italy, 2021, Revised Selected Papers (Lecture Notes in Business Information Processing, Vol. 436)*. Springer, 18–29.
- [36] Zhanke Zhou, Zhaocheng Zhu, Xuan Li, Mikhail Galkin, Xiao Feng, Sanmi Koyejo, Jian Tang, and Bo Han. 2025. Landscape of Thoughts: Visualizing the Reasoning Process of Large Language Models. arXiv:2503.22165 [cs.LG] <https://arxiv.org/abs/2503.22165>
- [37] Lisa Zimmermann, Francesca Zerbato, and Barbara Weber. 2023. What makes life for process mining analysts difficult? A reflection of challenges. *Software and Systems Modeling* (11 2023), 1–29.