# Evaluating Large Language Models in Process Mining: Capabilities, Benchmarks, and Evaluation Strategies

Alessandro Berti[1,2], Humam Kourani[1,2], Hannes Häfke[1], Chiao-Yun Li[1,2], Daniel Schuster[1,2]

[1] Fraunhofer FIT, Sankt Augustin, Germany
[2] Process and Data Science Chair, RWTH Aachen University, Aachen, Germany
{alessandro.berti,humam.kourani, hannes.haefke, chiao-yun.li, daniel.schuster}@fit.fraunhofer.de

**Abstract.** Using Large Language Models (LLMs) for Process Mining (PM) tasks is becoming increasingly essential, and initial approaches yield promising results. However, little attention has been given to developing strategies for evaluating and benchmarking the utility of incorporating LLMs into PM tasks. This paper reviews the current implementations of LLMs in PM and reflects on three different questions. 1) What is the minimal set of capabilities required for PM on LLMs? 2) Which benchmark strategies help choose optimal LLMs for PM? 3) How do we evaluate the output of LLMs on specific PM tasks? The answer to these questions is fundamental to the development of comprehensive process mining benchmarks on LLMs covering different tasks and implementation paradigms.

**Keywords:** Large Language Models (LLMs) · Output Evaluation · Benchmarking Strategies.

## 1   Introduction

Process mining (PM) is a data science field focusing on deriving insights about business process executions from event data recorded by information systems [1]. Several types of PM exist, including *process discovery* (learning process models from event data), *conformance checking* (comparing event data with process models), and *process enhancement* (adding frequency/performance metrics to process models). Although many automated methods exist for PM, human analysts usually handle process analysis due to the need for domain knowledge. Recently, LLMs have emerged as conversational interfaces trained on extensive data [28], achieving near-human performance in various general tasks [38]. Their potential in PM lies in *embedded domain knowledge* useful for generating database queries and insights [21], *logical and temporal reasoning capabilities* [2,16], *inference abilities over structured data* [12]. Prior research has asserted the usage of LLMs for PM tasks [3,4]. However, a comprehensive discussion on
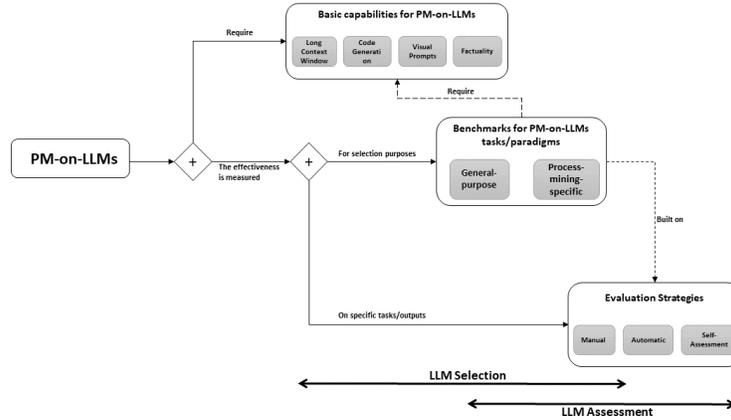
Fig. 1: Outline of the contributions of this paper.

necessary capabilities for PM, LLMs' suitability evaluation for process analytics, and assessment of LLMs' outputs in the PM context is lacking.

The three main contributions of this paper are summarized in Fig. 1. First, building upon prior work [3,4] proposing textual abstractions of process mining artifacts and an experimental evaluation of LLMs' responses, the essential capabilities that LLMs must have for PM tasks are derived in Section 3.1. The aforementioned capabilities allow us to narrow down the field of LLMs to those that meet these requirements. Next, evaluation benchmarks for selecting suitable LLMs are introduced in Section 3.2, incorporating both process-mining-specific and general criteria such as reasoning, visual understanding, factuality, and trustworthiness. Finally, we suggest automatic, human, and self-assessment methods for evaluating LLMs' outputs on specific tasks in Section 3.3, aiming to establish a comprehensive PM benchmark and enhance confidence in LLMs' usage, addressing potential issues like hallucination.

This paper provides an orientation to process mining researchers investigating the usage of LLMs, i.e., this paper aims to facilitate PM-on-LLMs research.

## 2   Background

LLMs enhance PM with superior capabilities, handling complex tasks through data understanding and natural language processing. This section covers PM tasks with LLMs (Section 2.1) and the adopted implementation paradigms (Section 2.2) along with the provision of additional domain knowledge.

### 2.1   Process Mining Tasks for LLMs

This subsection explores a range of PM tasks in which LLMs have already been adopted for process mining research. LLMs facilitate the automation of generating *textual descriptions* from process data, handling inputs such as event logs

or formal process models [4]. They also generate *process models* from textual descriptions, with studies showing LLMs creating BPMN models and declarative constraints from text [7]. In the realm of anomaly detection, LLMs play a crucial role in identifying process data *anomalies*, including unusual activities and performance bottlenecks, offering context-aware detection that adapts to new patterns through prompt engineering. This improves versatility over traditional methods [3,4]. For *root cause analysis*, LLMs analyze event logs to suggest causes of anomalies or inefficiencies, linking delays to specific conditions or events. This goes beyond predefined logic, employing language processing for context-aware analysis [3,4] In ensuring *fairness*, LLMs identify and mitigate bias in processes, suggesting adjustments. They analyze processes like recruitment to detect disparities in rejection rates or delays by gender or nationality, aiding in fair decision-making [22,4]. LLMs can also interpret and explain *visual data*, including complex visualizations, by describing event flows in dotted charts and identifying specific patterns, such as batch processing. For process improvement, after PM tasks identify and analyze problems, LLMs can suggest actions and propose new process constraints [22,4].

### 2.2   Implementation Paradigms of Process Mining on LLMs

To effectively employ LLMs for PM tasks, specific implementation paradigms are required [3,4]. This section outlines key approaches for implementing LLMs in PM tasks. We distinguish three main strategies:

- *Direct provision of insights*: A prompt is generated that merges data abstractions with a query about the task. Also, interactive dialogue between the LLM and the user is possible for step-by-step analysis. The user starts with a query and refines or adjusts it based on the LLM's feedback, continuing until achieving the desired detail or accuracy, such as pinpointing process inefficiencies. For instance, to have LLMs identifying unusual behavior in an event log, we combine a textual abstraction of the log (such as the directly-follows graph or list of process variants) with a question like "Can you analyze the log to detect any unusual behavior?"
- *Code generation*: LLMs can be used to create structured queries, like SQL, for advanced PM tasks [11]. Rather than directly asking LLMs for answers, users command LLMs to craft database queries from natural language. These queries are then executed on the databases holding PM information. It is applicable to PM tasks that can be converted into database queries, such as filtering event logs or computing the average duration of process steps. Also, LLMs can be used to generate executable programs that use existing PM libraries to infer insights over the event data [9].
- *Automated hypotheses generation*: Combining the previous strategies by using textual data abstraction to prompt LLMs for autonomous hypotheses generation [3,4]. The hypotheses are accompanied by SQL queries for verification against event data. Results confirm or refute these hypotheses, with potential for LLM-suggested refinements of hypotheses.

LLMs may require additional knowledge about processes and databases to implement PM tasks, for example, in anomaly detection and crafting accurate database queries. Some strategies are used to equip LLMs with this additional domain knowledge [14], including *fine-tuning* and *prompt engineering*.

## 3   Evaluating LLMs in Process Mining

This section introduces criteria for selecting LLMs that are suitable for PM tasks. Moreover, we introduce criteria for evaluating their outputs. First, in Section 3.1, we discuss the fundamental capabilities needed for PM (long context window, acceptance of visual prompts, coding, factuality). Then, we introduce in Section 3.2 general-purpose and process-mining-specific benchmarks to measure the different LLMs on process-mining-related tasks. To foster the development of process-mining-specific benchmarks and to be able to evaluate a given output, we propose in Section 3.3 different methods to evaluate the output of an LLM.

### 3.1   LLMs Capabilities Needed for Process Mining

In this section, we discuss four important capabilities of LLMs for PM tasks:

- *Long Context Window*: Event logs in PM often include a vast amount of cases and events, challenging the *context window* limit of LLMs, which restricts the token count in a prompt [13]. Moreover, also the textual specification of process models requires a significant amount of information. The context window limit can be severe in many currently popular LLMs.[3] Even simple abstractions like the ones introduced in [3] (directly-follows graph, list of process variants) may exceed this limitation. The context window, which is set during model training, must be large enough for the data size. Recent efforts aim to extend this limit, though quality may decline [13,20].
- *Accepting Visual Prompts*: Visualizations in PM, such as the dotted chart and the performance spectrum [15], summarizing process behavior, empower analysts to spot interesting patterns not seen in tables. Interpreting visual prompts is key for semi-automated PM. Large Visual Models (LVMs) use architectures similar to language models trained on annotated image datasets [31]. They perform tasks like object detection and image synthesis, recognizing patterns, textures, shapes, colors, and spatial relations.[4]
- *Coding (Text-to-SQL) Capabilities*: With the context window limit preventing full event log inclusion in prompts, generating scripts and database queries is crucial for analyzing event data. As discussed in Section 2.2, text-to-SQL assists in filtering and analyzing event data. Key requirements for text-to-SQL in PM include understanding database schemas, performing complex joins, using database-specific operators (e.g., for calculating date differences), and translating PM concepts into queries. Overall, modern LLMs offer excellent coding capabilities [3].

---

[3] `https://community.openai.com/t/are-the-full-8k-gpt-4-tokens-available-on-chatgpt/237999`

[4] GPT-4 and Google Bard/Gemini are popular models supporting both visual and textual prompts.

– *Factuality*: LLM hallucination involves generating incorrect or fabricated information [24]. Factuality measures an LLM's ability to cross-check its outputs against real facts or data, crucial for PM tasks like anomaly detection and root cause analysis. This may involve leveraging external databases [19], knowledge bases, or internet search [32] for validation. For instance, verifying the sequence Cancel Order" followed by Deliver Order" against public data in anomaly detection. LLMs with web browsing can access up-to-date information, enhancing factuality.[5]

### 3.2 Relevant LLMs Benchmarks

After identifying the required capabilities for LLMs in PM, benchmarking strategies are essential to measure the quality of the textual outputs returned by the LLMs satisfying such capabilities.

Considering the wide array of available benchmarks for assessing LLMs behavior, we focus on identifying those most relevant to PM capabilities. In [5], a comprehensive collection of benchmarks is introduced. This section aims to select and utilize some of these benchmarks to evaluate various aspects of LLMs' performance in PM contexts.

– *Traditional benchmarks*: Textual prompts are crucial for LLMs evaluation in PM. Benchmarks like AGIEval assess models via standardized exams [37], and MT-Bench focuses on conversational and instructional capabilities [36]. Another benchmark evaluates LLMs on prompts of long size [6].
– *Domain knowledge benchmarks*: Domain knowledge is essential for LLMs in PM to identify anomalies using metrics and context. Benchmarks like XIEZHI assess knowledge across different fields (economics, science, engineering) [8], while ARB evaluates expertise in areas like mathematics and natural sciences [26].
– *Visual benchmarks*: Understanding PM visualizations, such as dotted charts, is essential (c.f. Section 3.1). LLMs must accurately process queries on these visualizations. MMBench tests models on image tasks [17], and MM-Vet assesses recognition, OCR, among others [35]. Yet, they may not fully meet PM visualization analysis needs, particularly in evaluating line orientations and point size/color.
– *Benchmarks for Text-to-SQL*: In PM, generating SQL from natural language is key for tasks like event log filtering. Benchmarks such as SPIDER and SPIDER-realistic test LLMs on text-to-SQL conversion [23]. The APPS benchmark evaluates broader code generation abilities [10].
– *Fairness benchmarks*: they evaluate LLM fairness in PM by analyzing group treatment and bias detection. DecodingTrust measures LLM trustworthiness, covering toxicity, bias, robustness, privacy, ethics, and fairness [30].
– *Benchmarking the generation of hypotheses*: LLMs' ability to generate hypotheses from event data is vital to implement semi-autonomous PM agents.

---

[5] https://cointelegraph.com/news/chat-gpt-ai-openai-browse-internet-no-longer-limited-info-2021

Table 1: Implementation paradigms and benchmarks for LLMs in the context of different PM tasks.

| Task | Paradigms | | | Benchmarks Classes | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Direct Provision | Code Generation | Hypotheses Generation | Traditional | Domain Knowledge | Visual Prompts | Text-to-SQL | Fairness | Hypotheses Generation |
| Process Description | X | | | X | X | | | | |
| Process Modeling | X | X | X | X | X | | X | | X |
| Anomaly Detection | X | | X | X | X | | X | | X |
| Root Cause Analysis | X | | X | X | X | | X | | X |
| Ensuring Fairness | X | | X | X | X | | X | X | X |
| Expl. and Interpreting Visualizations | X | | | | X | X | | | |
| Process Improvement | X | X | X | X | X | | X | X | X |

While specific benchmarks for hypothesis generation are lacking, related studies like [29] and [34] evaluate LLMs using scientific papers.

In Table 1, we link process mining (PM) tasks to implementation paradigms and benchmarks. We discuss these tasks:

- *Process description* requires understanding technical terms relevant to the domain, crucial for accurately describing processes.
- *Process modeling* involves generating models from text, using SQL for declarative and BPMN XML for procedural models. LLMs should offer various model hypotheses.
- *Anomaly detection* and *root cause analysis* need domain knowledge to analyze process sequences or identify event attribute combinations causing issues.
- *Fairness* involves detecting biases by analyzing event attributes and values, necessitating hypothesis generation by LLMs.
- *Explaining and interpreting visualizations* requires extracting features from images and texts, offering contextual insights, like interpreting performance spectrum visualization [15].
- *Process improvement* entails suggesting text proposals or new constraints to enhance current models, leveraging code generation capabilities and understanding process limitations.

While general-purpose benchmarks are already developed and are easily accessible, they are not entirely suited for the task of PM-on-LLMs. In particular, visual capabilities (explaining and interpreting PM visualizations) and autonomous hypotheses generation require more PM-specific benchmarks. However, little research exists on PM-specific benchmarks [3,4].

### 3.3   How to Evaluate LLMs Outputs

This section outlines criteria for assessing the quality of outputs generated by LLMs in PM tasks, serving two primary objectives. The first objective is to as-

sist users in identifying and addressing hallucinations and inaccuracies in LLMs' outputs. The second aim is to establish criteria for developing an extensive benchmark specifically tailored to PM applications of LLMs. The strategies follow:

– *Automatic evaluation* is particularly suited for text-to-SQL tasks. In this context, the formal accuracy and conciseness (indicated by the length of the produced query) of the SQL queries generated can be efficiently assessed. Additionally, the creation of declarative constraints, designed to enhance process execution, can also be evaluated in terms of their formal correctness.
– *Human evaluation* is essential for LLM tasks like direct querying and hypothesis generation. For direct querying tasks such as anomaly detection and root cause analysis, important criteria are *recall* (the model's ability to identify expected insights) and *precision* (the correctness of insights). These criteria also apply to hypothesis generation. Additionally, evaluating the feedback cycle's effectiveness in validating original hypotheses is crucial for these tasks.
– *Self-evaluation* in LLMs tackles hallucinations, as noted by [24]. Techniques include *chain-of-thought*, where LLMs detail their reasoning, enhancing explanations [33]. *Confidence scores* let LLMs assess their insights' reliability, discarding uncertain outputs for quality [27]. *Ensembling*, or using results from multiple LLM sessions, increases accuracy via majority voting or confidence checks [18]. *Self-reflection*, an LLM reviewing its or another's output, detects errors [25]. In anomaly detection, using confidence scores to exclude doubtful anomalies and ensembling to confirm detections across sessions improves reliability.

## 4   Conclusion

This paper examines LLM applications in PM, offering three main contributions: identification of necessary LLM capabilities for PM, review of benchmarks from literature, and strategies for evaluating LLM outputs in PM tasks. These strategies aim to build confidence in LLM use and establish benchmarks to assess LLM effectiveness across PM implementations.

Our discussion centers on current generative AI capabilities within PM, anticipating advancements like deriving event logs from videos. Despite future enhancements, the criteria discussed here should remain pertinent. Benchmarking for PM tasks on large language models (LLMs) will evolve, including both general and PM-specific benchmarks, yet the foundational aspects and methodologies are expected to stay consistent.

## References

1. van der Aalst, W.M.P.: Process Mining - Data Science in Action, Second Edition. Springer (2016)

2. Bang, Y., Cahyawijaya, S., et al., N.L.: A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity (2023). https://doi.org/10.48550/arXiv.2302.04023
3. Berti, A., Qafari, M.S.: Leveraging Large Language Models (LLMs) for Process Mining (Technical Report) (2023). https://doi.org/10.48550/arXiv.2307.12701
4. Berti, A., Schuster, D., van der Aalst, W.M.P.: Abstractions, Scenarios, and Prompt Definitions for Process Mining with LLMs: A Case Study (2023). https://doi.org/10.48550/arXiv.2307.02194
5. Chang, Y., Wang, X., et al., J.W.: A Survey on Evaluation of Large Language Models (2023). https://doi.org/10.48550/arXiv.2307.03109
6. Dong, Z., Tang, T., et al., J.L.: BAMBOO: A Comprehensive Benchmark for Evaluating Long Text Modeling Capacities of Large Language Models (2023). https://doi.org/10.48550/arXiv.2309.13345
7. Grohs, M., Abb, L., et al., N.E.: Large Language Models Can Accomplish Business Process Management Tasks. In: BPM 2023 International Workshops. Lecture Notes in Business Information Processing, vol. 492, pp. 453–465. Springer (2023)
8. Gu, Z., Zhu, X., et al., H.Y.: Xiezhi: An Ever-Updating Benchmark for Holistic Domain Knowledge Evaluation (2023). https://doi.org/10.48550/arXiv.2306.05783
9. Härer, F.: Conceptual model interpreter for Large Language Models. In: ER 2023. CEUR Workshop Proceedings, vol. 3618. CEUR-WS.org (2023)
10. Hendrycks, D., Basart, S., et al., S.K.: Measuring Coding Challenge Competence With APPS. In: NeurIPS Datasets and Benchmarks 2021 (2021)
11. Jessen, U., Sroka, M., Fahland, D.: Chit-Chat or Deep Talk: Prompt Engineering for Process Mining (2023). https://doi.org/10.48550/arXiv.2307.09909
12. Jiang, J., Zhou, K., et al., Z.D.: StructGPT: A General Framework for Large Language Model to Reason over Structured Data. In: EMNLP 2023. pp. 9237–9251. Association for Computational Linguistics (2023)
13. Jin, H., Han, X., et al., J.Y.: LLM Maybe LongLM: Self-Extend LLM Context Window Without Tuning (2024). https://doi.org/10.48550/arXiv.2401.01325
14. Kampik, T., Warmuth, C., et al., A.R.: Large Process Models: Business Process Management in the Age of Generative AI (2023). https://doi.org/10.48550/arXiv.2309.00900
15. Klijn, E.L., Fahland, D.: Performance Mining for Batch Processing Using the Performance Spectrum. In: BPM 2019 International Workshops. Lecture Notes in Business Information Processing, vol. 362, pp. 172–185. Springer (2019)
16. Liu, H., Ning, R., et al., Z.T.: Evaluating the Logical Reasoning Ability of ChatGPT and GPT-4 (2023). https://doi.org/10.48550/arXiv.2304.03439
17. Liu, Y., Duan, H., et al., Y.Z.: MMBench: Is Your Multi-modal Model an All-around Player? (2023). https://doi.org/10.48550/arXiv.2307.06281
18. Lu, K., Yuan, H., et al., R.L.: Routing to the Expert: Efficient Reward-guided Ensemble of Large Language Models (2023). https://doi.org/10.48550/arXiv.2311.08692
19. Pan, S., Luo, L., et al., Y.W.: Unifying Large Language Models and Knowledge Graphs: A Roadmap (2023). https://doi.org/10.48550/arXiv.2306.08302
20. Peng, B., Quesnelle, J., et al., H.F.: YaRN: Efficient Context Window Extension of Large Language Models (2023). https://doi.org/10.48550/arXiv.2309.00071
21. Petroni, F., Rocktäschel, T., et al., S.R.: Language Models as Knowledge Bases? In: EMNLP-IJCNLP 2019. pp. 2463–2473. Association for Computational Linguistics (2019)
22. Qafari, M.S., van der Aalst, W.M.P.: Fairness-Aware Process Mining. In: CoopIS 2019. Lecture Notes in Computer Science, vol. 11877, pp. 182–192. Springer (2019)

23. Rajkumar, N., Li, R., Bahdanau, D.: Evaluating the Text-to-SQL Capabilities of Large Language Models (2022). https://doi.org/10.48550/arXiv.2204.00498
24. Rawte, V., Chakraborty, S., et al., A.P.: The Troubling Emergence of Hallucination in Large Language Models - An Extensive Definition, Quantification, and Prescriptive Remediations. In: EMNLP 2023. pp. 2541–2573. Association for Computational Linguistics (2023)
25. Ren, J., Zhao, Y., et al., T.V.: Self-Evaluation Improves Selective Generation in Large Language Models (2023)
26. Sawada, T., Paleka, D., et al., A.H.: ARB: Advanced Reasoning Benchmark for Large Language Models (2023). https://doi.org/10.48550/arXiv.2307.13692
27. Singh, A.K., Devkota, S., et al., B.L.: The Confidence-Competence Gap in Large Language Models: A Cognitive Study (2023). https://doi.org/10.48550/arXiv.2309.16145
28. Teubner, T., Flath, C.M., et al., C.W.: Welcome to the Era of ChatGPT et al. Bus. Inf. Syst. Eng. **65**(2), 95–101 (2023)
29. Tong, S., Mao, K., Huang, Z., Zhao, Y., Peng, K.: Automating Psychological Hypothesis Generation with AI: Large Language Models Meet Causal Graph (2023)
30. Wang, B., Chen, W., et al., H.P.: DecodingTrust: A Comprehensive Assessment of Trustworthiness in GPT Models (2023). https://doi.org/10.48550/arXiv.2306.11698
31. Wang, J., Liu, Z., et al., L.Z.: Review of Large Vision Models and Visual Prompt Engineering (2023). https://doi.org/10.48550/arXiv.2307.00855
32. Wang, L., Ma, C., et al., X.F.: A Survey on Large Language Model based Autonomous Agents (2023). https://doi.org/10.48550/arXiv.2308.11432
33. Wei, J., Wang, X., et al., D.S.: Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In: NeurIPS 2022 (2022)
34. Yang, Z., Du, X., et al., J.L.: Large Language Models for Automated Open-domain Scientific Hypotheses Discovery (2023). https://doi.org/10.48550/arXiv.2309.02726
35. Yu, W., Yang, Z., et al., L.L.: MM-Vet: Evaluating Large Multimodal Models for Integrated Capabilities (2023). https://doi.org/10.48550/arXiv.2308.02490
36. Zheng, L., Chiang, W., et al., Y.S.: Judging LLM-as-a-judge with MT-Bench and Chatbot Arena (2023). https://doi.org/10.48550/arXiv.2306.05685
37. Zhong, W., Cui, R., et al., Y.G.: AGIEval: A Human-Centric Benchmark for Evaluating Foundation Models (2023). https://doi.org/10.48550/arXiv.2304.06364
38. Zhou, Y., Muresanu, A.I., et al., Z.H.: Large Language Models are Human-Level Prompt Engineers. In: ICLR 2023. OpenReview.net (2023)