

# On LLM-Assisted Generation of Smart Contracts from Business Processes

Fabian Stiehle<sup>1</sup>, Hans Weytjens<sup>1</sup>, and Ingo Weber<sup>1,2</sup>

<sup>1</sup> Technical University of Munich, School of CIT, Germany, [first.last@tum.de](mailto:first.last@tum.de)

<sup>2</sup> Fraunhofer Gesellschaft, Munich, Germany

**Abstract.** Large language models (LLMs) have changed the reality of how software is produced. Within the wider software engineering community, among many other purposes, they are explored for code generation use cases from different types of input. In this work, we present an exploratory study to investigate the use of LLMs for generating smart contract code from business process descriptions, an idea that has emerged in recent literature to overcome the limitations of traditional rule-based code generation approaches. However, current LLM-based work evaluates generated code on small samples, relying on manual inspection, or testing whether code compiles but ignoring correct execution. With this work, we introduce an automated evaluation framework and provide empirical data from larger data sets of process models. We test LLMs of different types and sizes in their capabilities of achieving important properties of process execution, including enforcing process flow, resource allocation, and data-based conditions. Our results show that LLM performance falls short of the perfect reliability required for smart contract development. We suggest future work to explore responsible LLM integrations in existing tools for code generation to ensure more reliable output. Our benchmarking framework can serve as a foundation for developing and evaluating such integrations.

**Keywords:** Blockchain, Process Execution, Process Enactment, Workflow, Large language models, Generative AI

## 1 Introduction

In many fields, LLMs are envisioned to assist in a wide variety of tasks, where so far the training of general machine learning models has been challenging due to the scarcity of specialized data and the large computational effort required (c.f., [4,8,42]).

In the wider field of software engineering, LLMs are anticipated to support all phases of the software engineering process [4]. For code generation, LLMs have arguably made the biggest practical impact so far, testified by the integration of commercial tools like *Github’s Copilot*—which builds on seminal research on LLMs trained on code [11]—into popular development environments like *Visual Studio Code*.

Similarly, within the field of business process management (BPM), Vidgof et al. [42] call to evaluate the combination of LLMs with existing BPM technologies. First results evaluating LLMs on BPM tasks show the large promise, since they perform comparable to or better than existing BPM tools [17].

Blockchain-based business process execution relies on a model-driven paradigm, where process descriptions are transformed into executable artefacts based on rule-based transformation tools [39]. These tools, however, exhibit various limitations such as in their flexibility, e.g., in terms of supported process modelling constructs, their supported output targets, or their support of blockchain-specific features. Considering the wider field of model-driven engineering, LLMs are similarly prospected to have the potential to drive automation [8]. This hope can draw on seminal results from the related field of code-to-code translation (transpilers), where LLM-based approaches were found to outperform traditional rule-based approaches [34,35].

While many positive visions exist, and early results on leveraging LLMs to assist in code generation are impressive [7], significant challenges remain—extending even beyond the well-known hallucination issue (or more precisely: confabulation [36]). For instance, GitHub Copilot can introduce numerous security vulnerabilities into generated code [32]. LLM outputs are inherently non-deterministic [30], making them unreliable for consistent behaviour. Meanwhile, Huang et al. [20] show that generated code may reproduce ethically concerning biases, such as gender-related ones.

Furthermore, proprietary models raise concerns about confidentiality, privacy, and autonomy (c.f. [4,18]). These AI models are often deployed on large, centralised platforms provided by hyperscalers like *AWS*, *Azure*, or *GCP*. Open-source models running on these platforms face the same security concerns. Relying on central deployments may not be a good fit for blockchain-based processes, where decentralisation is a goal [38].

Thus, it is paramount to systematically evaluate LLMs’ usefulness and fit of properties for a given task and weigh benefits and drawbacks. Research is called to identify the scenarios where LLMs add true value (c.f. [4,42]). However, evaluation poses a significant challenge, as it often requires manual human investigation and large volumes of labelled or parallel data [9].

Research on using LLMs for smart contract generation is in its infancy. Existing work aims, among other things, to reduce the required specialised programming skills for contract development. However, most of this work focuses on reporting the syntactical correctness of generated code (compilability), not its functional correctness [1,25,28]. In other cases, correctness is assessed through manual inspection of small samples [10,15]. In the case of blockchain-based process execution, early work explores the promise of deriving code from process descriptions, but draws conclusions from a singular case study [15]. Current research furthermore does not make data or code openly accessible.

In this work, we present an exploratory study investigating the use of LLMs to generate smart contract code from business process descriptions. We think it paramount to ground such research in open empirical evidence from larger datasets. Beyond standards of open science, the availability of open repeatable benchmarks and data is especially important in the present case, as LLMs capabilities develop in a fast pace and their output is unpredictable. With this work, we introduce an automated evaluation framework for generating smart contract code from process models. We provide empirical data from larger data sets of process models (165 models filtered and sampled from the collection of SAP-SAM [37] models). We test seven LLMs

of different types and sizes in their capabilities of achieving central properties of blockchain-based process enactment (c.f. [39]): enforcing process flow, case data-based conditions, resource allocation, and efficiency.

In general, our results indicate good performance of some models, with those achieving F1 scores of 0.8 or more. Due to the stochastic nature of LLMs, output remains imperfect and unreliable. While such performance may be acceptable in other contexts, blockchain is an unforgiving environment—on public blockchains developers should assume that any weakness will be exploited. As such, we believe this is a fundamental issue of the chosen approach to have LLMs generate smart contract code, which cannot be overcome by LLMs based on the current architectures. We discuss this point further in the paper, including an outlook on roles LLMs could fulfil productively in the generation of smart contracts from process models.

The remainder of the paper is structured as follows. A background on LLMs and relevant terminology is provided in the next section; note that we assume familiarity with blockchain and process enactment on it. Subsequently related work is discussed in Section 3, before we present the benchmarking framework in Section 4. Based on it, we conduct experiments that we report on in Section 5 and discuss in Section 6, before Section 7 concludes.

## 2 Background

In this section, we give relevant background on the AI models, their attributes, and related concepts which we use in the body of the paper. Large Language Models (LLMs) are transformer-based neural network systems [41]. LLMs are a specific class of Foundation Models, a class of machine learning models trained on extensive data sets, not for one specific purpose but as a basis for many possible applications [3]. In the case of LLMs, training ingests very large textual corpora, comprising not only natural language, but also programming code (such as Solidity for blockchain smart contracts), and formal representations (such as BPMN for modeling business processes). Due to their generality and the training input, LLMs can perform complex tasks beyond language understanding and generation, including analyzing and creating business process models [19] and smart contracts [12]. Many of the latest LLMs support other modalities (such as images, video, or audio), hence they are sometimes referred to as large multi-modal models (LMMs). However, for this paper, the distinction between LLMs and LMMs is of no importance, and hence for the sake of clear communication, we follow the current common practice and refer to them as LLMs.

In many usage scenarios, LLMs deliver high performance out of the box, eliminating the need for fine-tuning or training from scratch with large supervised datasets or significant compute resources. In *zero-shot prompting*, models complete tasks based solely on instructions without prior examples, whereas *few-shot prompting* involves providing a handful of illustrative examples directly within the input [14].

*Tokens* are the units of input and output for LLMs. Tokens can be whole words, subwords (parts of words), individual characters, punctuation, or special characters [2]. The total number of tokens—both in the input query and the generated output—directly influences the computational load of running these models. Models

that perform complex *reasoning* (so called reasoning models) typically generate a plan to answer a query, execute the steps in the plan, and possibly check their work; hence they require many more tokens than regular LLMs. Increasing the model size—measured by the number of parameters, the numeric values representing the strength of connections between “neurons”—generally improves performance but also raises computational demands, leading to greater energy consumption [29]. Balancing model size, token usage, and capability is therefore essential for developing efficient and sustainable LLM-based applications.

LLMs are available in both proprietary and open-source forms [43]. Proprietary models, such as *OpenAI’s GPT* and *Anthropic’s Claude* model families, are typically accessed via APIs hosted by third parties. In comparison to open-source models, they often deliver superior performance but introduce risks such as data exposure, dependency on external providers, and limited transparency. In contrast, open-source models enable self-hosting and on-premise deployment—an attractive option in blockchain contexts where data sovereignty, trust minimization, and operational independence are essential.

*Temperature* is a setting for LLMs that controls the randomness of generated text; lower temperatures make outputs more predictable and focused (by selecting the most probable tokens), while higher temperatures encourage more creative and varied responses (by selecting less probable tokens).

### 3 Related Work

Within the field of BPM, Vidgof et al. [42] outline the opportunities and challenges of integrating LLM-based tools within the BPM lifecycle. Recent work explores an increasingly wide array of BPM applications (see e.g., Pfeiffer et al. [33], which explore four real-world use cases that demonstrate the use of LLMs across modelling, prediction and automation). Within predictive process monitoring, LLMs are explored for their capability to predict future states of processes (e.g., Pasquadibisceglie et al. [31]). In prescriptive process monitoring, LLMs are explored to enhance recommendations with LLM-generated explanations [23].

Orthogonal to our work is the question whether LLMs can assist in deriving accurate models from natural language process descriptions (e.g., Hörner et al [19]); for an overview see Klietsova et al. [22]. Closer to our work, Monti et al. [26] propose an LLM-driven pipeline to extract executable scripts (for deployment in a process execution engine) from natural language process descriptions. For the evaluation of their code generation, they use 10 cases and compare the LLM output to a manually implemented script. Additionally, they conducted a human evaluation, assessing the quality of the produced code. A similar study to ours was conducted by Berti et al. [5], which presents a benchmark on the performance of LLMs for different Process Mining tasks, the open-ended nature of these tasks requires a *LLM judge* evaluation approach, where LLMs assess the output of other LLMs.

Research on LLM-assisted smart contract generation is in its infancy. Existing results are limited to reporting the syntactical correctness of generated code (compilability) and the automated detection of known vulnerabilities using existing

tools [1,28,25]. A relevant result to our investigation is provided by Luo et al. [25]; their result suggests that augmenting the generation process by a formalised model improves results. Some works go beyond syntactical correctness by manual inspection of code [10,15]. However, this is limited to a few simplistic cases, as it is highly labour intensive. Karanjai et al. [21] extract solidity functions from GitHub repositories and use extracted code comments as prompts; to test the LLM output, they rely on corresponding unit tests present in the mined repository.

None of the aforementioned studies on smart contract generation make their data or evaluation frameworks available.

Our work lies at the intersection of previously outlined fields. To the best of our knowledge, the only directly related work is the recent case study of Gao et al. [15]. They present an approach, based on few-shot prompting, capable of generating a smart contract from a Business Process Modeling Language (BPML) specification (a superset of BPEL), derived from a collaboration diagram. However, their evaluation is limited to one process, which includes start, end, and message events, tasks, and two XOR splits. They do not make their code or data available.

To the best of our knowledge, we are the first to present an open source benchmarking framework to automate the assessment of code generation from process models. We provide an instantiation of it based on the large SAP-SAM data set [37], one-shot and two-shot prompts, and a Ethereum virtual machine (EVM) environment. Using it, we evaluate different proprietary and open source LLMs on 165 cases each. All our data is available and our tests can be repeated (see Footnote 3).

## 4 Benchmarking Framework

To assess the capabilities of LLMs for smart contract generation, we designed a configurable benchmarking framework and make it openly available along with all relevant input and output data.<sup>3</sup> An open framework facilitates repeatability and can be used to judge the capabilities of future model evolutions. We envision our framework to also serve as a foundation for evaluation in future related research (as outlined in Section 6). In contrast to previous work, we want to facilitate automated evaluation from large data sets with optimal coverage. An established method to benchmark the correctness of a blockchain-based business process is to replay all possible conforming traces<sup>4</sup> (which the smart contract has to accept) and replay a set of non-conforming traces (which the smart contract has to reject) [39]. Our approach makes use of this method. From a given process model, it generates conforming traces (not always exhaustively, but up to a configurable threshold, as parallelism and loops

<sup>3</sup> <https://github.com/fstiehle/bpmn-sol-llm-benchmark>. An archived version is available at: <https://doi.org/10.5281/zenodo.16616694>.

<sup>4</sup> For the remainder of the paper, we use common terminology. We loosely denote an event log as a set of events, where each event is associated with a case identifier that groups it into a case. A trace is the ordered sequence of events (activities) that occurred for a specific case. Each event represents a task (activity) in the model. A trace can be said to be in conformance with a process model if it represents a valid execution path through that model.

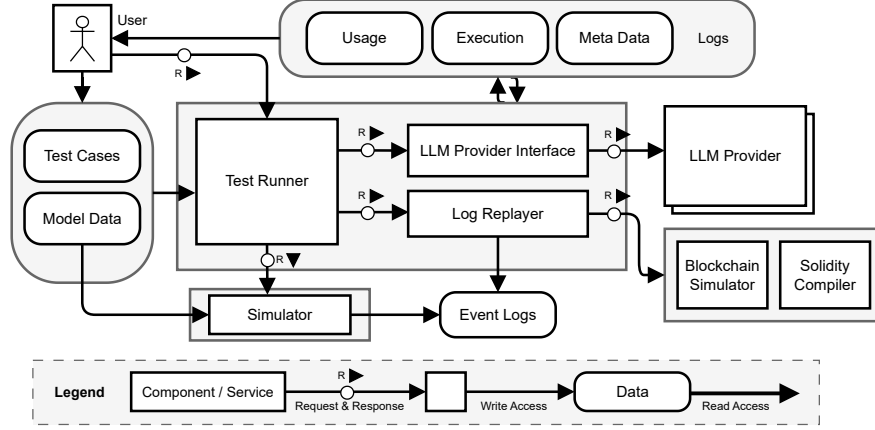


Fig. 1: Main components, services and data of the benchmarking framework architecture as FMC block diagram.

can result in an intractable search space) and non-conforming traces. All traces are then replayed against the generated smart contracts.

#### 4.1 Architecture Overview

Figure 1 depicts an overview of our architecture. As input, different test cases can be configured. Mainly, a test case is a tuple of: (i) the LLM to benchmark, (ii) the prompt to use, (iii) the process model data set to use.

Interaction with the framework is achieved through a *test runner*. Based on user interaction, the runner coordinates the benchmark execution with the other internal and external components. From the model data (process models), a *simulator* component generates conforming and non-conforming traces. This simulator component is external to the framework, so it can be swapped based on different model data inputs. The simulator is responsible for generating conforming and non-conforming event logs from the input process models. To generate a non-conforming trace, a conforming trace is randomly chosen and manipulated. The simulator also generates an encoding that maps the events and participants to how they should be represented in the smart contract (`taskIDs` and `participantIDs`, the latter associated with a blockchain address).

This encoding is embedded into the prompt, along with the model data, by the *LLM provider interface*, which calls the external LLM provider as configured in the test cases. The received output is stored in a usage log. For any interaction with an LLM provider, the usage log stores (among other things): a timestamp, the configured test case, the full input and the full output. This allows to reconstruct and rerun a benchmark at a later time. The usage log also records usage statistics (tokens used, price, etc.) as reported by the provider.

From the usage log output, the *log replayer* extracts the smart contract code, compiles and deploys it to an external blockchain environment, and uses the encodings

and event logs to perform a benchmark on the deployed contract. The replayer stores the results in execution logs. The framework also extracts metadata on the used process model (e.g, modelling constructs per model, etc.) to aid in the evaluation of the benchmark.

## 4.2 Instantiation

In the current instantiation of our framework, we support BPMN 2.0 Choreographies. This is a purely practical implementation choice, given that there is no consensus on the best fitting modelling paradigm for blockchain-based execution (c.f. [39]) and given our familiarity with a suitable tool. We instantiate our framework for an Ethereum virtual machine (EVM) blockchain environment, the most widely employed environment [39]. We discuss key design decisions in the instantiation of our architecture.

*Simulator.* We extend the open source tool *Chorpiller*, first introduced in [40] with simulation capabilities. Chorpiller transforms BPMN Choreographies to smart contracts, generates non-conforming traces from conforming traces, and also generates machine-readable encodings on how to interact with a contract. Chorpiller parses a given Choreography into an interaction net, a special type of labelled Petri net (see [13]), suitable to represent choreographies. We make use of this intermediate presentation to generate conforming traces. Here, we adopt the implementation of *pm4py* [6], a popular Python library for process mining, which includes a *playout* functionality to generate event log traces from Petri nets. The implementation tries to discover new conforming traces with each pass until a threshold of passes is reached.

As we also want to benchmark data-based exclusive gateways (XOR), we had to extend the playout functionality to generate appropriate data manipulation events. To do so, for each outgoing flow (other than the default flow), we generate a boolean decision. During trace generation, once our algorithm encounters a decision transition, it inserts a corresponding data event in the trace.<sup>5</sup>

*Test Runner, Replayer and LLM Provider.* To provide the replayer with a blockchain environment, we use *hardhat*, a popular Ethereum development framework that allows testing, deployment, and debugging of smart contracts in a locally simulated EVM environment.<sup>6</sup> As LLM provider, we use *OpenRouter*, a platform that provides a unified API across multiple language model providers. This simplifies our integration and gives access to a broad range of state-of-the-art models.<sup>7</sup> We use a *Node.js* environment; the test runner is implemented using *Mocha*, a JavaScript test framework.<sup>8</sup>

<sup>5</sup> Chorpiller implements transactional logic as in [24]; the smart contract makes as much progress as possible after a task is executed; i.e., data-based decisions are made autonomously once the gateway is enabled, and the required data must be already set by then. Thus, data events are inserted preceding any event that leads to a given gateway.

<sup>6</sup> <https://hardhat.org>, accessed 2025-06-12

<sup>7</sup> <https://openrouter.ai>, accessed 2025-06-12

<sup>8</sup> <https://www.npmjs.com>, <https://mochajs.org>, both accessed 2025-06-12

## 5 Experiment

We use the instantiation of our framework to conduct a large scale benchmarking experiment. Our process model data is based on the *SAP Signavio Academic Models Dataset* (SAP-SAM), which was initially introduced in [37]. The dataset contains different model types created through the *Academic Initiative* platform from 2011 to 2021.<sup>9</sup> Thus, the dataset primarily contains process models created by students, researchers, and teaching staff. Still, for BPMN specifically, the dataset’s properties (distribution of modelling constructs) are in line with previous research assessing the usage of modelling constructs from diverse sources [27]. The collection includes 4,096 BPMN 2.0 choreography models, to our knowledge, the largest collection of choreography models accessible for research purposes.

### 5.1 Pre-Processing

The SAP-SAM dataset contains many non-standard compliant choreography models. For our purposes, we can relax requirements of ownership and observability present in the standard, as the smart contract provides global ownership and observability (c.f. [24]).<sup>10</sup> Furthermore, the models lack execution-relevant information for exclusive gateways (conditions or labels from which conditions could be inferred, and default flow markings). Thus, we pre-processed each model. When no default flow was marked, we set the first outgoing flow to the default flow. Then, for all other outgoing flows, we inserted a boolean condition. Furthermore, we removed any Signavio extension elements. To reduce the size of LLM input, we also removed the BPMN 2.0 Diagram Interchange, as it only contains additional information required to visualise the model. Finally, we merged all start and end events, so each model contains only one, to adhere to the implementation limitation of Chorpiller.

From these pre-processed models, we use Chorpiller to assess the syntactic soundness of each model. Simply put, if Chorpiller is able to generate a contract from the model, we consider it for our benchmark. Chorpiller supports all basic elements of BPMN Choreographies<sup>11</sup> Choreography tasks, start and end event, exclusive, event, and parallel gateways, sub choreographies, and loops in the model. It also ignores issues regarding ownership and observability.

After the filtering steps, 1,427 choreography models remain. We use a sample of 165 models for our benchmark runs. In our sample, on average, each process contains 13 participants (see footnote 10), six tasks, one diverging exclusive gateway, 0.1 diverging event-based gateways, and 0.2 diverging parallel gateways, among others. The largest model in the dataset contains 24 tasks and ten gateways, respectively.

<sup>9</sup> <https://academic.signavio.com>, accessed 2025-06-12

<sup>10</sup> During our initial exploration of the dataset we encountered this issue in many models. The problem is exacerbated by the fact that many participants share the same name in the task bands, but are assigned different participant IDs. We hypothesise that this occurred when a participant’s name was entered manually rather than selected from a list of existing ones in the visual editor.

<sup>11</sup> We use the latest alpha version (<https://github.com/fstiehle/chorpiller/tree/release/v2>, accessed 2025-06-12).



	Provider	Model	Size (B)
Open Source	DeepSeek	DeepSeek-V3 0324	671
	Meta	Llama-3.1-405b-instruct	405
		Llama-3.3-70b-instruct	70
	Alibaba	Qwen3-235b-a22b	235
Proprietary	OpenAI	GPT-4.1	n/a
	Anthropic	Claude Sonnet 4	n/a
	X AI	Grok 3	n/a

Table 1: Models selected for our experiments, conducted June 11–13, 2025, using the then-current models available through OpenRouter. *Size* refers to the number of parameters in billions (B), where available.

## 5.2 Benchmark

*Prompt.* For our benchmark, we developed multiple prompts. To reduce the likelihood of our results being tainted by a poorly performing prompt, we tested, compared, and refined multiple versions in pre-runs, which we conducted with sets of five to twenty process models. This allowed us to iteratively refine our prompts and test our framework. For any pre-run, in addition to the automated tests, we manually investigated the generated output. This led to many refinements. In summary, our prompts moved from loosely defined (zero-shot) instructions, to (better performing) more specific instructions on how the process model should be interpreted and the contract generated. Through our initial tests, we arrived at a one-shot prompt. Specifically, in our prompt we ask for a Solidity implementation for the given process model, enforcing: (i) the control flow, i.e., the order of tasks, (ii) that only the respective initiator can execute a task, and (iii) the autonomous enforcement of gateways, and the evaluation of data-based decisions.<sup>12</sup> To gauge the effect of prompt-based training, we test a one-shot and a two-shot variant.

*Setup.* We select a range of top proprietary LLMs<sup>13</sup> as benchmarks and compare them to a host of open source models to evaluate to what extent hosting open source models signifies sacrificing performance for autonomy. Table 1 provides an overview of our benchmarked models. Figure 2 gives an overview of our benchmarking experiment run. For our LLM selection (c.f. Table 1), we benchmark a one-shot and two-shot variant of our prompt with temperature set to 0,<sup>14</sup> on a sample of 165 process models from our pre-processed files. For the generation of conforming process traces, we set a

<sup>12</sup> Our prompt also specifies that the state of the contract should be encoded using a bitmasking technique, as it is the most efficient encoding for a token-based execution (c.f. [16]). This variant did, on average, not perform worse than a prompt asking for a more naïve implementation during our pre-runs.

<sup>13</sup> Our initial run also included *Google’s Gemini*. However, we were not able to deactivate the output of its reasoning process, which prevented us from reliably parsing a contract from the output automatically.

<sup>14</sup> Leading to quasi-deterministic (c.f. [30]) inference results. Tie breaking between tokens of equal probabilities, floating point variability, etc. may still cause stochasticity.

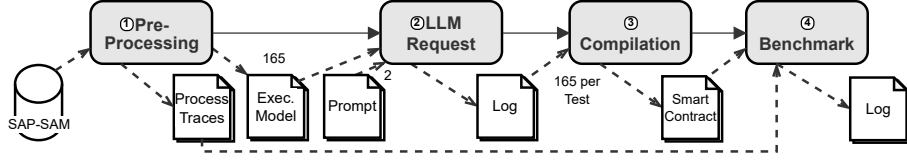


Fig. 2: Process of our benchmark experiment: We pre-process the raw model data from the SAP-SAM dataset to receive executable process models and corresponding traces. We use 165 models and perform requests for one and two-shot prompts, a pair for each LLM. From this, logs are generated including the received output. All output is compiled and deployed and benchmarked against the process traces. The benchmark result is logged.

threshold of 2,500 traces per process. We generated and replayed 50 non-conforming traces per process.<sup>15</sup>

We ran our experiment from June 11 to June 13, 2025. For some requests, we had to perform repeated tries, as the provider connection sometimes timed out. We also encountered a period in which OpenRouter experienced an outage. Our framework is set up to retry failed requests on additional runs.

### 5.3 Results

To assess the quality of the produced output we compare usage, in terms of cost and tokens used, as charged by OpenRouter.<sup>16</sup> To assess the efficiency of the generated output, we also record gas usage. For a given process model and its generated implementation, we classify the outcome of a trace replay accordingly.

- *True Positive*: Each event in a conforming trace was accepted (led to a state change in the contract), and the whole trace led to the end event.
- *False Positive*: A non-conforming trace was accepted as per above.
- *True Negative*: Any event in the non-conforming trace was rejected, or the trace did not lead to the end event.
- *False Negative*: A conforming trace was rejected as per above.

Using this classification framework, we calculate precision and recall per process case, using standard formulations, and the F1 score (the harmonic mean of precision

<sup>15</sup> For process models with loops, the threshold is a necessary upper bound. Our non-conforming trace generation algorithm currently depends on a ground truth to assess whether a manipulated trace is not conforming on accident. For models, where the number of conforming traces exceeds our search threshold, non-conforming traces that are actually conforming can be generated. These must be manually removed.

<sup>16</sup> OpenRouter operates on a ‘credits’ system; credits are purchased upfront, which are then deducted per model request, according to the underlying provider’s token-based rate. OpenRouter charges a fee ( $\approx 5\% + \$0.35$ ) when loading credits., which we did not include in our calculations (<https://openrouter.ai/docs/faq>, accessed 2025-06-13).

Shot	Model	(Avg./Process)		Correctness	
		Cost(\$)	Tokens	F1 Mac.	Comp.(%)
One	grok-3-beta	0.044	10.134	0.918	100.0
	claude-sonnet-4	0.046	11.442	0.862	100.0
	gpt-4.1	0.028	10.326	0.797	99.4
	qwen3-235b-a22b	0.009	18.444	0.648	97.0
	deepseek-chat-v3-0324	0.005	10.483	0.580	99.4
	llama-3.1-405b-instruct	0.010	10.259	0.475	99.4
	llama-3.3-70b-instruct	0.001	10.249	0.399	90.4
Two	grok-3-beta	0.056	14.964	0.861	100.0
	claude-sonnet-4	0.064	17.218	0.853	100.0
	gpt-4.1	0.038	15.410	0.696	100.0
	deepseek-chat-v3-0324	0.007	15.680	0.669	97.6
	qwen3-235b-a22b	0.009	24.415	0.581	93.4
	llama-3.1-405b-instruct	0.016	15.473	0.431	98.8
	llama-3.3-70b-instruct	0.002	15.252	0.370	97.0

Table 2: Result of our benchmark run with 165 process models. We report the average cost in US-\$ (as reported by OpenRouter), and tokens used per process model. The overall correctness is reported via the F1 macro. Compilability (Comp.) reports on the percentage of syntactically correct generated contracts.

and recall), common metrics to assess LLM output (see e.g., [9]). In our case, F1 is a suitable choice over metrics like accuracy, since the number of traces is unbalanced.

To assess the overall performance of each model, we calculate the macro F1 (the average of all F1 across all cases). Our results are shown in Table 2. The most obvious aspect is that Grok and Claude achieved F1 scores of 0.8 or more in all variants. Most of the generated code compiled. The cost for translating a process model was on the orders of 0.1 cents to a few cents. Interestingly, the two-shot prompt did not consistently yield better results. As a side effect of the experiment, we observed that the framework performed well in running the benchmark across the diverse set of process and AI models.

## 6 Discussion, Limitations & Future Work

Our results show that current LLMs can transform executable choreography models into syntactically and functionally correct smart contracts most of the time—even when benchmarked against a realistic and diverse dataset. While the results in terms of F1 score shows promise, it can only serve as a point of departure for future work exploring meaningful integrations within a smart contract generation workflow. As mentioned earlier, blockchain is an unforgiving environment, and smart contract vulnerabilities may become very costly. Indeed, the goal of business process execution via smart contracts is to provide a trusted and secure decentralised platform. These requirements would be undermined by a reckless integration of LLM capabilities.

F1 scores that do not reliably achieve 100% would not be suitable for this context. Say, the approaches would be improved to achieving an average F1 score of 98%; while this would be impressive in many domains, it falls short of the perfect reliability required for blockchain-based smart contracts. Given the financial risks

and immutable nature of blockchain transactions, even such a 2% error rate could lead to significant vulnerabilities or losses, making such performance inadequate for real-world deployment. We do not see a way in which this fundamental issue could be resolved with current LLM architectures.

However, we see ways to advance the current direction beyond using output as-is. Future work should explore integrating LLM capabilities into existing smart contract generation tools. For smart contract development, which demands high security, LLM integration must rely on extensive evaluation and robust verification of generated outcomes. This could involve using LLMs to propose code snippets or modifications, which are then rigorously checked against formal specifications (as we demonstrated with our framework) or verified using automated theorem provers, before being considered. LLMs should also generate test cases or identify potential vulnerabilities themselves, identifying common or context-depending issues in smart contracts, augmenting existing verification processes.

Furthermore, LLMs could be used to extend the functionality of existing code generation tools (generating new rules) by: (i) generating more flexible templates based on specific process models, (ii) generating code snippets for edge cases not covered by standard templates, (iii) suggesting optimisations for rules and generated code, and (iv) assisting in translating between different smart contract languages or blockchain platforms. Any such use would still have to be vetted in a suitable form, but has the advantage that rule improvements and extensions are not subject to non-determinism after being included in the code generation tools. This approach combines the reliability and domain-specific knowledge of traditional code generation tools with the flexibility and natural language understanding of LLMs, and hence addresses the limitations mentioned in the introduction.

Our benchmarking framework can serve as a foundation for evaluating these future directions, helping to assess the quality of LLM-generated code, the effectiveness of verification methods, and benchmarking extended code generation tools. Towards this, the capabilities of the benchmarking tool itself must be extended to, e.g., consider other factors such as efficiency of generated code, potential biases present in the LLM, or sustainability factors (c.f. [18]).

Finally, some limitations apply to our performed benchmark. We experimented with different prompts, but cannot guarantee that the selected query was optimal; furthermore, we used the same prompt across all LLMs, which may be suboptimal. Although we included a diverse set of LLMs, our findings should not be assumed to generalise to all current or future LLMs.

## 7 Conclusions

In this work, we presented an exploratory study investigating the use of LLMs for generating smart contract code from business process descriptions. We introduced an automated evaluation framework and provided empirical data from a large dataset of 165 process models. Our results show that while current LLMs can transform executable choreography models into syntactically and functionally correct smart contracts most of the time, achieving F1 scores of 0.8 or more for top-performing models, this performance falls short of the perfect reliability required for smart

contracts. Given the financial risks and immutable nature of blockchain transactions, even small error rates could lead to significant vulnerabilities or losses. We argue that this fundamental issue cannot be resolved with current LLM architectures. Instead, we propose future work to explore responsible LLM integrations in existing tools for code generation, focusing on using LLMs for verification and enhancing current code generation tools rather than replacing them entirely. Our benchmarking framework can serve as a foundation for developing and evaluating such integrations.

**Acknowledgments.** Generative AI was used to assist in the editing of the manuscript and the implementation of the artifact. Generated output was never taken "as-is", it was reviewed and verified by the authors.

## References

1. Alam, M.T., Goswami, S., Singh, K., Halder, R., Maiti, A., Banerjee, S.: Solgen: Secure smart contract code generation using large language models via masked prompting. In: Proceedings of the 18th Innovations in Software Engineering Conference. pp. 1–11 (2025)
2. Ali, M., Fromm, M., Thellmann, K., et al.: Tokenizer choice for LLM training: Negligible or crucial? In: Findings of the Association for Computational Linguistics: NAACL 2024. pp. 3907–3924 (2024)
3. Bass, L., Lu, Q., Weber, I., Zhu, L.: Engineering AI Systems: Architecture and DevOps Essentials. Addison-Wesley Professional (2025)
4. Belzner, L., Gabor, T., Wirsing, M.: Large language model assisted software engineering: prospects, challenges, and a case study. In: International Conference on Bridging the Gap between AI and Reality. pp. 355–374. Springer (2023)
5. Berti, A., Kourani, H., van der Aalst, W.M.: PM-LLM-Benchmark: Evaluating large language models on process mining tasks. In: International Conference on Process Mining. pp. 610–623. Springer (2024)
6. Berti, A., van Zelst, S., Schuster, D.: PM4Py: A process mining library for python. *Software Impacts* **17**, 100556 (2023)
7. Brynjolfsson, E., Li, D., Raymond, L.R.: The productivity effects of generative AI: Evidence from a field experiment with GitHub Copilot (2023), <https://mit-genai.pubpub.org/pub/v5iixkxsv/release/2>, accessed 2025-06-07
8. Burgueño, L., Di Ruscio, D., Sahraoui, H., Wimmer, M.: Automation in model-driven engineering: A look back, and ahead. *ACM Transactions on Software Engineering and Methodology* (2025)
9. Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., Chen, H., Yi, X., Wang, C., Wang, Y., et al.: A survey on evaluation of large language models. *ACM transactions on intelligent systems and technology* **15**(3), 1–45 (2024)
10. Chatterjee, S., Ramamurthy, B.: Efficacy of various large language models in generating smart contracts. In: Future of Information and Communication Conference. pp. 482–500. Springer (2025)
11. Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H.P.D.O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al.: Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021)
12. De Vito, G., D’Amici, D., Izzo, F., Ferrucci, F., Di Nucci, D.: LLM-based generation of solidity smart contracts from system requirements in natural language: The AstraKode case. In: International Conference on Software Analysis, Evolution and Reengineering (SANER). pp. 170–180. IEEE (2025)

13. Decker, G., Weske, M.: Local enforceability in interaction Petri nets. In: BPM. LNCS, vol. 4714, pp. 305–319. Springer, Cham (2007)
14. Dong, Q., Li, L., Dai, D., Zheng, C., Ma, J., Li, R., Xia, H., Xu, J., Wu, Z., Liu, T., Chang, B., Sun, X., Li, L., Sui, Z.: A survey on in-context learning (2024), <https://arxiv.org/abs/2301.00234>
15. Gao, S., Liu, W., Zhu, J., Dong, X., Dong, J.: BPMN-LLM: Transforming bpmn models into smart contracts using large language models. *IEEE Software* (2025)
16. García-Bañuelos, L., Ponomarev, A., Dumas, M., Weber, I.: Optimized Execution of Business Processes on Blockchain. In: BPM, pp. 130–146. Springer, Cham (2017)
17. Grohs, M., Abb, L., Elsayed, N., Rehse, J.R.: Large language models can accomplish business process management tasks. In: International Conference on Business Process Management. pp. 453–465. Springer (2023)
18. Haase, J., Klessascheck, F., Mendling, J., Pokutta, S.: Sustainability via LLM right-sizing. *arXiv preprint arXiv:2504.13217* (2025)
19. Hörner, L.F.: Towards an LLM-based conversational framework for business process modeling: Research approach and preliminary results. In: International Conference on Advanced Information Systems Engineering. pp. 286–293. Springer (2025)
20. Huang, D., Zhang, J.M., Bu, Q., Xie, X., Chen, J., Cui, H.: Bias testing and mitigation in LLM-based code generation. *ACM Transactions on Software Engineering and Methodology* (2024)
21. Karanjai, R., Li, E., Xu, L., Shi, W.: Who is smarter? an empirical study of ai-based smart contract creation. In: 2023 5th Conference on Blockchain Research & Applications for Innovative Networks and Services (BRAINS). pp. 1–8. IEEE (2023)
22. Klievtsova, N., Benzin, J.V., Kampik, T., Mangler, J., Rinderle-Ma, S.: Conversational process modelling: state of the art, applications, and implications in practice. In: International Conference on Business Process Management. pp. 319–336. Springer (2023)
23. Kubrak, K., Botchorishvili, L., Milani, F., Nolte, A., Dumas, M.: Explanatory capabilities of large language models in prescriptive process monitoring. In: Business Process Management. pp. 403–420. Springer Nature Switzerland, Cham (2024)
24. Ladleif, J., Weske, M., Weber, I.: Modeling and Enforcing Blockchain-Based Choreographies. In: BPM, vol. 11675, pp. 69–85 (2019)
25. Luo, H., Lin, Y., Yan, X., Hu, X., Wang, Y., Zeng, Q., Wang, H., Jiang, J.: Guiding LLM-based smart contract generation with finite state machine. *arXiv preprint arXiv:2505.08542* (2025)
26. Monti, F., Leotta, F., Mangler, J., Mecella, M., Rinderle-Ma, S.: Nl2processops: towards LLM-guided code generation for process execution. In: BPM. pp. 127–143 (2024)
27. Muehlen, M.z., Recker, J.: How Much Language Is Enough? Theoretical and Practical Use of the Business Process Modeling Notation, pp. 429–443. Springer (2013)
28. Napoli, E.A., Barbàra, F., Gatteschi, V., Schifanella, C.: Leveraging large language models for automatic smart contract generation. In: Computers, Software, and Applications Conference (COMPSAC). pp. 701–710. IEEE (2024)
29. O'Donnell, J., Crownhart, C.: We did the math on AI's energy footprint. here's the story you haven't heard. *MIT Technology Review* (May 2025)
30. Ouyang, S., Zhang, J.M., Harman, M., Wang, M.: An empirical study of the non-determinism of chatgpt in code generation. *ACM Transactions on Software Engineering and Methodology* **34**(2), 1–28 (2025)
31. Pasquadibisceglie, V., Appice, A., Malerba, D.: LUPIN: A LLM approach for activity suffix prediction in business process event logs. In: ICPM. pp. 1–8 (2024)
32. Pearce, H., Ahmad, B., Tan, B., Dolan-Gavitt, B., Karri, R.: Asleep at the keyboard? Assessing the security of GitHub Copilot's code contributions. *Communications of the ACM* **68**(2), 96–105 (2025)

33. Pfeiffer, P., Rombach, A., Majlatow, M., Mehdiyev, N.: From theory to practice: Real-world use cases on trustworthy LLM-driven process modeling, prediction and automation (2025), <https://arxiv.org/abs/2506.03801>
34. Roziere, B., Lachaux, M.A., Chausson, L., Lample, G.: Unsupervised translation of programming languages. *Advances in neural information processing systems* **33**, 20601–20611 (2020)
35. Roziere, B., Zhang, J.M., Charton, F., Harman, M., Synnaeve, G., Lample, G.: Leveraging automated unit tests for unsupervised code translation. *arXiv preprint arXiv:2110.06773* (2021)
36. Smith, A.L., Greaves, F., Panch, T.: Hallucination or confabulation? Neuroanatomy as metaphor in large language models. *PLOS Digital Health* **2**(11), 1–3 (11 2023)
37. Sola, D., Warmuth, C., Schäfer, B., Badakhshan, P., Rehse, J., Kampik, T.: SAP Signavio academic models: A large process model dataset. In: *Process Mining Workshops at ICPM 2022. LNBIP*, vol. 468, pp. 453–465. Springer (2022)
38. Stiehle, F., Klessascheck, F., Kjær, M., Weber, I.: Business in the age of platform economics: Managing decentralised business processes beyond blockchain. In: *Innov8BPM'24: International Workshop on Managing Process Innovation and Value Creation in the Era of Digital Transformation at BPM'24* (2024)
39. Stiehle, F., Weber, I.: Blockchain for business process enactment: a taxonomy and systematic literature review. In: *BPM: Forum. LNBIP*, vol. 459, pp. 5–20 (2022)
40. Stiehle, F., Weber, I.: Process channels: A new layer for process enactment based on blockchain state channels. In: *BPM*. pp. 198–215. Springer (2023)
41. Vaswani, A., Shazeer, N., Parmar, N., et al.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
42. Vidgof, M., Bachhofner, S., Mendling, J.: Large language models for business process management: Opportunities and challenges. In: *International Conference on Business Process Management*. pp. 107–123. Springer (2023)
43. Wong, E.: Comparative analysis of open source and proprietary large language models: Performance and accessibility. *Advances in Computer Sciences* **7**(1), 1–7 (2024)