

DABL: Detecting Semantic Anomalies in Business Processes Using Large Language Models

Wei Guan¹, Jian Cao^{1*}, Jianqi Gao¹, Haiyan Zhao², Shiyu Qian¹

¹ Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China

² Department of Computer Science and Engineering, University of Shanghai for Science and Technology, Shanghai, China
{guan-wei, cao-jian, 193139, qshiyu}@sjtu.edu.cn, zhaohaiyan1992@foxmail.com

Abstract

Detecting anomalies in business processes is crucial for ensuring operational success. While many existing methods rely on statistical frequency to detect anomalies, it's important to note that infrequent behavior doesn't necessarily imply undesirability. To address this challenge, detecting anomalies from a semantic viewpoint proves to be a more effective approach. However, current semantic anomaly detection methods treat a trace (i.e., process instance) as multiple event pairs, disrupting long-distance dependencies. In this paper, we introduce DABL, a novel approach for detecting semantic anomalies in business processes using large language models (LLMs). We collect 143,137 real-world process models from various domains. By generating normal traces through the playout of these process models and simulating both ordering and exclusion anomalies, we fine-tune Llama 2 using the resulting log. Through extensive experiments, we demonstrate that DABL surpasses existing state-of-the-art semantic anomaly detection methods in terms of both generalization ability and learning of given processes. Users can directly apply DABL to detect semantic anomalies in their own datasets without the need for additional training. Furthermore, DABL offers the ability to interpret anomalies' causes in natural language, providing valuable insights into the detected anomalies.

Code — <https://github.com/guanwei49/DABL>

Introduction

Business process anomaly detection is geared towards identifying undesired behavior occurring during process execution, serving as a crucial component in guaranteeing the efficient and dependable operation of businesses. By pinpointing anomalies within business processes, these detection techniques facilitate timely intervention, maintenance, and optimization, consequently bolstering overall well-being (Guan et al. 2024d).

Over the past few decades, notable advancements have been achieved in business process anomaly detection. Fig. 1 illustrates a comparison of various applicable methods for accomplishing this task. Traditional statistical-based approaches (Lu, Fang, and Fang 2022; Ko and Comuzzi 2022;

Nolle et al. 2022; Guan et al. 2024a) rely on analyzing statistical frequencies to identify anomalies. However, infrequent behavior is not necessarily anomalous; it may represent rare but acceptable behavior. Conversely, frequent behavior may not always be normal. Furthermore, these methods focus on providing anomaly scores and require manual specification of thresholds to distinguish between normal and anomalous instances, which is not suitable for real-world applications. Alternatively, methods based on conformance checking (Ebrahim and Golpayegani 2022; Sarno, Sinaga, and Sungkono 2020; Sinaga and Sarno 2016) detect anomalies by assessing the alignment between traces and their corresponding process models. Yet, accurately capturing complex processes within a process model remains a challenge, thereby restricting the utility of such approaches. The concept of **semantic anomaly detection**, recently introduced, addresses these challenges by identifying anomalies from a semantic viewpoint. For example, it can detect irregularities such as a claim being paid after having been rejected. Its grounding in natural language analysis allows for the consideration of typical behavior in standard processes, eliminating the necessity of having a specific process model at hand. However, existing semantic-based anomaly detection methods (van der Aa, Rebmann, and Leopold 2021; Caspary, Rebmann, and van der Aa 2023) treat a trace as multiple event pairs, disrupting long-distance dependencies and thus limiting their accuracy. Additionally, these methods interpret the cause of anomalies by providing anomalous event pairs, which can be confusing.

Recently, there have been significant advancements in LLMs. Due to their remarkable language comprehension abilities, LLMs such as GPT-3.5 (Ouyang et al. 2022), GPT-4 (Achiam et al. 2023), Llama 2 (Touvron et al. 2023), and GLM-3 (Zeng et al. 2022) have shown proficiency in tasks like summarization, paraphrasing, and instruction following in zero-shot scenarios. However, in the context of semantic anomaly detection in business processes, their performance is limited by a lack of prior domain knowledge. As illustrated in Fig. 1, they often struggle to provide specific answers.

To address the aforementioned issues, we propose DABL, a fine-tuned LLM designed to detect semantic anomalies in business processes. Due to the lack of event logs comprising traces from various domains with rich semantic in-

*Corresponding author.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

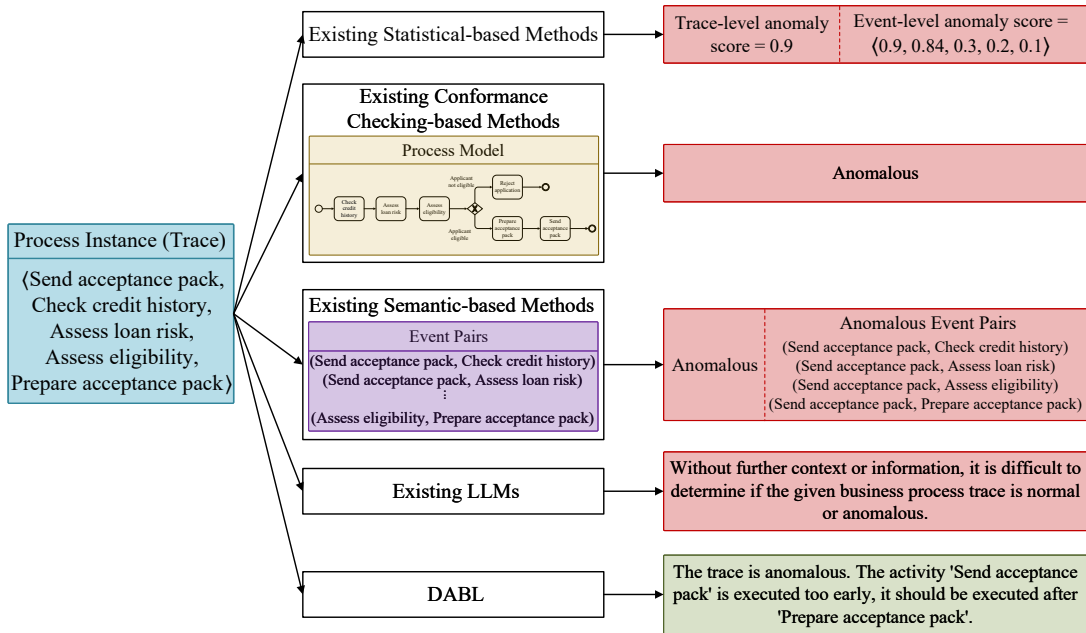


Figure 1: Comparison between our DABL with existing methods.

formation, we generated our training dataset by playing out 143,137 real-world process models from three different process model datasets. This resulted in 1,574,381 normal traces. The collected process models cover a broad range of domains, including common processes related to order and request handling, as well as specialized processes in fields such as software engineering and healthcare. Utilizing the generated normal traces, we then created synthetic anomalous traces. We introduced ordering anomalies, where activities should be executed in a different sequence (e.g., "accept request" followed by "check request"), and exclusion anomalies, where certain activities should not occur together within the same trace without an intermediate activity (e.g., "refusing the application" followed by "accepting the application" without "reapplying" in between). These generated normal and anomalous traces collectively form the training dataset. Finally, by incorporating traces into question and answer content, we fine-tune the Llama 2-Chat 13B model (Touvron et al. 2023), an open-source LLM, using QLoRA (Dettmers et al. 2024), to create a generic model capable of detecting semantic anomalies in business processes. Compared to existing anomaly detection methods, DABL offers the capability to interpret the causes of anomalies in natural language, providing valuable insights into the detected anomalies. Extensive experiments show that DABL surpasses state-of-the-art methods in both generalization ability and learning of given processes.

Thanks to DABL's strong generalization ability, users can apply our open-source, fine-tuned model directly to their datasets without the need for additional training. Notably, it operates in a zero-shot manner, meaning it does not need normal traces or a process model during its operation. Our contributions are summarized as follows:

- We introduce DABL, an innovative method for fine-

tuning LLMs to detect semantic anomalies in business processes.

- We introduce novel techniques for simulating business process anomalies, encompassing ordering anomalies and exclusion anomalies, thereby enabling precise fine-tuning of LLMs.
- Extensive experiments show that DABL outperforms state-of-the-art methods in both generalization ability and the learning of given processes. Tests on real-world datasets confirm the practical effectiveness of DABL.

Related Work

Business Process Anomaly Detection

Existing business process anomaly detection methods can be divided into three categories: statistical-based, conformance checking-based, and semantic-based.

Statistical-based Methods Some of these methods construct probabilistic models to infer the probability values (i.e., anomaly score) of traces. For example, HPDTMC (Yang et al. 2020) constructs discrete-time Markov chains (DTMC) and introduces hitting probabilities (HP). EDBN (Pauwels and Calders 2019a,b) extends dynamic Bayesian networks by adding functional dependencies. PN-BBN (Lu, Fang, and Fang 2022) extends Petri nets by incorporating Bayesian networks. Other methods convert traces into vector representations and then detect anomalies using data mining techniques such as local outlier factor (LOF) and isolation forest (IF). For example, activities are considered as words and encoded into vectors using the word2vec (Mikolov et al. 2013) technique in (Junior et al. 2020; Vertuam Neto et al. 2021). Additionally, Trace2vec (De Koninck, vanden Broucke, and De Weerd 2018), extended from

doc2vec and n-gram, is used to encode traces in (Rullo et al. 2020). Furthermore, the authors in (Ko and Comuzzi 2022, 2021) employ one-hot encoding to convert traces into vector representations and detect anomalies using statistical leverage (Hoaglin and Welsch 1978). Recently, deep learning has been adopted to detect anomalies based on reconstruction errors. Given that traces exhibit sequential data characteristics, the authors in (Guan et al. 2023; Nolle et al. 2022; Krajsic and Franczyk 2021; Guan et al. 2024c) embed LSTM, GRU or transformers within autoencoders to enhance the model’s reconstruction capabilities. The authors in (Huo et al. 2021; Guan et al. 2024a; Niro and Werner 2024) transform traces into graphs and utilize graph neural networks (GNNs) to generate graph encodings, identifying anomalies by evaluating the reconstruction error of the graphs or traces.

Statistical-based methods detect anomalies by analyzing statistical frequencies. However, infrequent behavior is not necessarily anomalous, as it may represent rare but acceptable behavior. Conversely, frequent behavior may not always be normal.

Conformance Checking-based Methods The conformance checking-based approaches (Ebrahim and Golpayegani 2022; Sarno, Sinaga, and Sungkono 2020; Sinaga and Sarno 2016) utilize process models, which are either provided by the user or derived from logs using process mining techniques. Anomalies are detected through conformance checking techniques (Leemans, Fahland, and van der Aalst 2018), which assess the alignment between traces and the corresponding process model. When the trace deviates from the process model, it is considered anomalous.

The performance of these methods heavily depends on the quality of the process model. Additionally, complex processes are difficult to accurately represent with a process model, limiting the applicability of these methods.

Semantic-based Methods The semantic-based methods detect anomalies through natural language analysis, aiming to detect process behaviors that deviate from a semantic point of view.

Van der Aa et al. (van der Aa, Rebmann, and Leopold 2021) fine-tune BERT (Devlin et al. 2018), a pre-trained language model, to parse the names of executed activities by extracting the action and business object. Then, a knowledge base capturing assertions about the interrelations that should hold among actions parsed from names of executed activities is applied. The knowledge is extracted either from VerbOcean (Chklovski and Pantel 2004) or from an abstract representation of the process model. Anomalies can be detected by checking if the recorded process behavior violates the assertions captured in the knowledge base.

Caspary et al. (Caspary, Rebmann, and van der Aa 2023) extract event pairs that are in an eventually-follow relation. To detect anomalous event pairs, they propose two approaches: an SVM-based approach and a BERT-based approach. The SVM-based approach transforms an event pair into a vector representation using GloVe embeddings. This vector is then fed into a trained SVM, which classifies whether the event pair is an anomaly. The BERT-based approach extends BERT with an additional output layer for

two-class classification, determining whether an input event pair is anomalous or not. Both the SVM and the extended BERT model are trained using normal event pairs extracted from normal traces, along with anomalous event pairs simulated by randomly generating event pairs that are not normal.

However, existing semantic-based methods treat a trace as multiple event pairs, which disrupts long-distance dependencies. Additionally, these methods only identify anomalous event pairs to interpret the causes of anomalies, making them difficult to understand. In contrast, our DABL incorporates the entire trace into a novel prompt, allowing the LLMs to capture long-distance dependencies. DABL also provides insightful interpretations of the causes of anomalies in natural language, making them easy to understand.

Large Language Models for Anomaly Detection

Motivated by the impressive cognitive abilities exhibited by LLMs (Ouyang et al. 2022; Zeng et al. 2022; Achiam et al. 2023; Touvron et al. 2023), researchers have begun investigating their application for anomaly detection. AnomalyGPT (Gu et al. 2024) and Myriad (Li et al. 2023) incorporate novel image encoders with LLMs for industrial anomaly detection (IAD). Elhafsi et al. (Elhafsi et al. 2023) apply an LLM to analyze potential confusion among observed objects in a scene, which could lead to task-relevant errors in policy implementation. LLMAD (Liu et al. 2024) leverages LLMs for few-shot anomaly detection by retrieving and utilizing both positive and negative similar time series segments. In (Qi et al. 2023; Guan et al. 2024b; Egersdoerfer, Zhang, and Dai 2023), authors explore the application of LLMs for log-based anomaly detection. SheepDog (Wu and Hooi 2023) conducts fake news detection by preprocessing data using LLMs to reframe the news, customizing each article to match different writing styles. Sarda et al. (Sarda et al. 2023) propose a pipeline for automatic microservice anomaly detection and remediation based on LLMs.

Yet, the application of LLMs for business process anomaly detection remains unexplored.

Method

DABL is a novel conversational fine-tuned large language model, primarily designed to detect semantic anomalies in business processes and interpret their causes. Fig. 2 details the DABL training procedure, which consists mainly of dataset preparation and fine-tuning.

Dataset Preparation

To effectively fine-tune LLMs for developing a generic model capable of detecting semantic anomalies in business processes, a log meeting the following criteria is imperative: i) it must encompass both normal and anomalous traces, ii) it should contain rich semantic information (i.e., the activities should not be represented by meaningless characters), and iii) the traces within it should stem from diverse processes across various domains. Since such a log is not available in the real world, we generate normal traces by playout of the real-world process models from the BPM Academic Initiative (BPMAI) (Weske et al. 2020), fundamentals of business

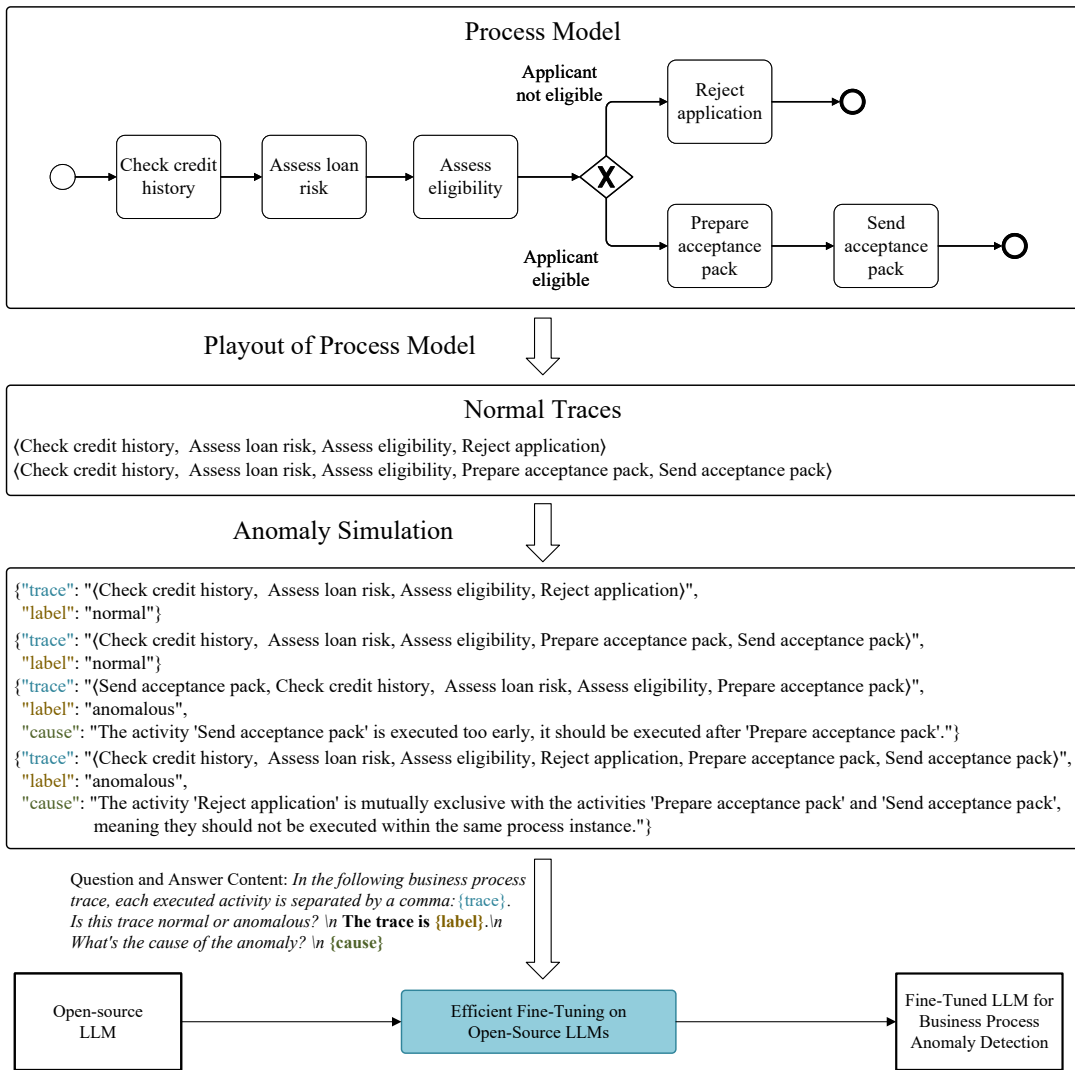


Figure 2: Detecting semantic anomalies in business processes using large language models.

process management (FBPM) (Dumas et al. 2018), and SAP signavio academic models (SAP-SAM) (Sola et al. 2022). These process models cover a broad range of domains, including common processes related to order and request handling, as well as specialized processes from fields such as software engineering and healthcare. We then generate synthetic anomalies from these normal traces. We detail the dataset preparation in the following subsection.

Generation of Normal Traces We select process models \mathcal{M} from BPMAI, FBPM, and SAP-SAM that meet the following criteria: they are in BPMN notation (Chinosi and Trombetta 2012), described in English, and convertible into a sound workflow net (Van Der Aalst et al. 2011). This results in a total set of 144,137 process models. Among these, 143,137 process models are used for generating training datasets, resulting in 1,574,381 normal traces, while the remaining 1,000 process models are used for generating test datasets.

Next, for each process model $m \in \mathcal{M}$, we perform a playout (Berti, Van Zelst, and van der Aalst 2019) to obtain the set of normal traces, denoted as \mathcal{L}_m . These traces are allowed by the process model m . To prevent infinite trace lengths, we limit each loop in the process model m to be executed a maximum of twice.

Anomaly Simulation Normal traces can be converted into anomalous ones by disrupting the order of executed activities (ordering anomalies) or by sequencing exclusive activities (exclusion anomalies).

Ordering anomalies: Ordering anomalies arise when activities ought to be executed in a different sequence. Five types of ordering anomalies, as identified in (Nolle et al. 2022), are frequently encountered in real-world business processes. Fig. 3 illustrates distinct anomalous traces resulting from the application of these five anomaly types to a normal trace $\langle e_1, e_2, \dots, e_n \rangle$. These anomaly types are defined as follows and are applied to a normal trace to generate

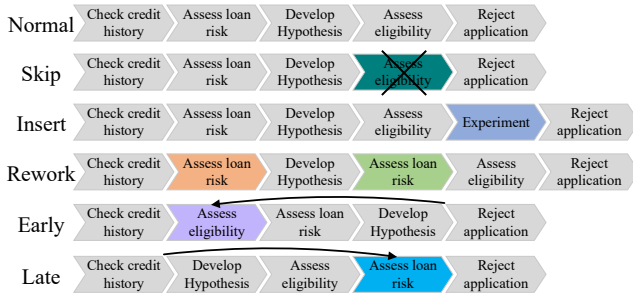


Figure 3: Different anomaly types applied to a normal trace.

ordering anomalies:

- *Skip*: A sequence of up to three activities $\langle e_i, \dots, e_j \rangle$ is skipped.
- *Insert*: A sequence of up to three random activities $\langle e'_1, \dots, e'_m \rangle$ is inserted. The random activities are selected from a set comprising all possible activities across all process models.
- *Rework*: A sequence of up to three activities $\langle e_i, \dots, e_j \rangle$ is executed a second time after activity e_k .
- *Early*: A sequence of up to three activities $\langle e_i, \dots, e_j \rangle$ is executed earlier and consequently skipped later.
- *Late*: A sequence of up to three activities $\langle e_i, \dots, e_j \rangle$ is executed later and consequently skipped earlier.

Below are the causes for these anomaly types, which are currently in the plural form (i.e., 'activities', 'they'). During implementation, they may need to be flexibly transformed into the singular form (i.e., 'activity', 'it').

- *Skip*: The activities $\{e_i, \dots, e_j\}$ are skipped before $\{e_{j+1}\}$.
- *Insert*: The activities $\{e'_1, \dots, e'_m\}$ should not be executed.
- *Rework*: The activities $\{e_i, \dots, e_j\}$ are reworked after $\{e_k\}$.
- *Early*: The activities $\{e_i, \dots, e_j\}$ are executed too early, they should be executed after $\{e_{i-1}\}$.
- *Late*: The activities $\{e_i, \dots, e_j\}$ are executed too late, they should be executed before $\{e_{j+1}\}$.

Here, $\{e_i, \dots, e_j\}$ represents converting the trace $\langle e_i, \dots, e_j \rangle$ into a string format by enclosing each executed activity in apostrophes and separating them with commas, while using *and* before the penultimate and ultimate activities. For example, for the trace $\langle A, B, C, D \rangle$, the resulting string would be 'A', 'B', 'C' and 'D'.

However, the generated anomalies may actually represent a normal trace. For a process model m , we refine the set of generated ordering anomalies \mathcal{L}_m^o by excluding traces present in \mathcal{L}_m from it.

Exclusion anomalies: Exclusion anomalies occur when certain activities should not have been executed together within the same trace without an intermediate activity. For instance, in the loan application process illustrated in Fig. 2, it is inappropriate to send an acceptance pack and reject an application within the same trace.

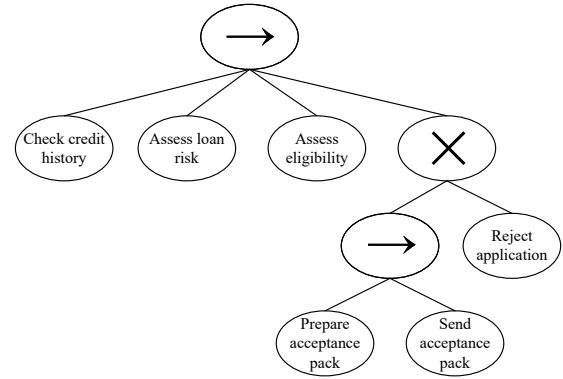


Figure 4: Different anomaly types applied to a normal trace.

The process tree (Aalst, Buijs, and Dongen 2011), a specialized form of process model, is utilized for analyzing process structure. For instance, the process tree corresponding to the process model depicted in Fig. 2 is illustrated in Fig. 4. We begin by converting the gathered process model m into a process tree using the techniques presented in (van Zelst and Leemans 2020). Then, we replace an exclusive node (represented as \times) in the process tree with a parallel node (represented as \wedge), resulting in a modified process model denoted as m' . This modification enables certain exclusive activities to be executed within the same trace to simulate exclusion anomalies. It is important to recognize that a single process tree may contain multiple exclusive nodes; therefore, we carry out this modification successively, resulting in multiple modified process models. Subsequently, we ploy out of all the modified models, restricting each loop to be executed a maximum of twice, to generate the set of traces, denoted as $\mathcal{L}_{m'}$. We refine $\mathcal{L}_{m'}$ by excluding traces present in \mathcal{L}_m from it, resulting in the set of exclusion anomalies \mathcal{L}_m^e .

To extract the causes of exclusion anomalies, we need to identify activities that exhibit exclusion relationships. In a process tree, activities located under different branches of an exclusive node (represented as \times) exhibit such relationships. For example, in the process tree illustrated in Fig. 4, the activity set $\{\text{Prepare acceptance pack}, \text{Send acceptance pack}\}$ and the activity set $\{\text{Reject application}\}$ exhibit exclusion relationships. Formally, consider a modified model m' , which results from modifying an exclusive node R in the process tree corresponding to the process model m . An exclusion anomaly t is generated from model m' . The node R has N branches, with the activity sets $\mathcal{A}_1, \dots, \mathcal{A}_N$ under them. Activities within each activity set \mathcal{A}_i that do not appear in t are filtered out. The cause of this exclusion anomaly t is then:

- The activities $\{\mathcal{A}_1\}$ are mutually exclusive with the activities $\{\bigcup_{i=2, \dots, N} \mathcal{A}_i\}$, meaning they should not be executed within the same process instance.

Question and Answer Content To conduct prompt tuning on the LLM, we generate corresponding textual queries based on simulated anomalous traces. Specifically, each query consists of two components.

The first component introduces the traces, such as "In the

following business process trace, each executed activity is separated by a comma: *(Send acceptance pack, Check credit history, Assess loan risk, Assess eligibility, Prepare acceptance pack)*". The second component queries whether the trace is anomalous, asking, for instance, "Is this trace normal or anomalous?". The LLM first responds to whether the given trace is normal or anomalous. If it is anomalous, the LLM is asked about the cause of the anomaly, for example, "What causes this trace to deviate?". The LLM then interprets the cause of the anomaly, such as "The activity 'Send acceptance pack' is executed too early, it should be executed after 'Prepare acceptance pack'.". This content about the anomaly's cause provides valuable insights and facilitates actions to maintain the health of the process execution.

Efficient Fine-Tuning on LLMs

DABL involves fine-tuning the open-source Llama 2-Chat 13B model (Touvron et al. 2023) to enhance its capability to detect semantic anomalies in business processes. To mitigate the expense associated with fine-tuning LLMs with a substantial parameter count, we leverage QLoRA (Dettmers et al. 2024) to reduce memory usage. QLoRA achieves this by back-propagating gradients into a frozen 4-bit quantized model while preserving the performance level attained during the full 16-bit fine-tuning process.

We employ the Adam optimizer (Kingma and Ba 2014) to fine-tune the LLMs for two epochs, setting the initial learning rate to 5×10^{-5} with polynomial learning rate decay. The mini-batch size is set to 64. The fine-tuning is carried out on an NVIDIA A6000 GPU with 48 GB of memory.

Experiments

Experimental Setup

Datasets As mentioned in the previous section, we allocate 1,000 process models for generating the test dataset \mathcal{D}_1 . These models produce 14,387 normal traces, and we randomly simulate anomalies, resulting in 13,694 anomalous traces. In total, the test dataset \mathcal{D}_1 comprises 28,081 traces.

From 143,137 process models used for generating the training dataset, we randomly select 1,000 process models to create the test dataset \mathcal{D}_2 . These 1,000 process models produce 21,298 normal traces, and we randomly simulate anomalies, resulting in 19,627 anomalous traces. In total, the test dataset \mathcal{D}_2 comprises 40,925 traces. Note that, although the normal traces within the test dataset \mathcal{D}_2 are identical to those in the training dataset, the simulated anomalies are not.

In summary, the test dataset \mathcal{D}_1 is used to evaluate the model's generalization ability, verifying if the model can detect anomalies of unseen processes. The test dataset \mathcal{D}_2 aims to validate the model's performance on seen processes but unseen anomalies (i.e., learning of given processes).

Compared Methods Statistical-based and conformance checking methods can only be applied to datasets containing traces from a single process. However, our test datasets include traces from 1000 processes where no two traces are identical (i.e., traces with identical orders of activities are executed). Therefore, these methods cannot be compared.

	Prec.(%)	Rec.(%)	F ₁ (%)	Acc.(%)
SEM	48.67	46.8	47.72	50.81
SENSE-SVM	87.95	1.12	2.20	52.50
SENSE-BERT	48.17	97.74	64.53	48.47
DBAL	94.06	89.79	91.88	92.39

Table 1: Semantic anomaly detection results on dataset \mathcal{D}_1 . The best results are indicated using bold typeface.

	Prec.(%)	Rec.(%)	F ₁ (%)	Acc.(%)
SEM	71.91	48.63	58.02	66.75
SENSE-SVM	90.28	28.64	43.49	64.82
SENSE-BERT	93.16	62.88	75.08	80.28
DBAL	98.12	95.64	96.87	97.03

Table 2: Semantic anomaly detection results on dataset \mathcal{D}_2 . The best results are indicated using bold typeface.

In our evaluation, we compare our DABL to existing semantic business process anomaly detection methods: SENSE (Caspary, Rebmann, and van der Aa 2023) and SEM (van der Aa, Rebmann, and Leopold 2021). SENSE offers both SVM-based and BERT-based models for detecting anomalous event pairs, which we denote as SENSE-SVM and SENSE-BERT, respectively. These methods divide traces into event pairs and determine whether each pair is normal or anomalous. If at least one event pair in a trace is identified as anomalous, the entire trace is classified as anomalous. It is important to note that SEM can only detect anomalous event pairs that share the same business object, automatically classifying pairs with distinct business objects as non-anomalous. Due to the high training costs, we utilize the open-source trained models provided by the authors for dataset \mathcal{D}_1 . For dataset \mathcal{D}_2 , we train the comparative models using the 21,298 normal traces available within it. The hyper-parameters of these methods are set to the values that yielded the best results reported in the original paper.


Evaluation Metrics Following existing anomaly detection methods, we employ *precision*, *recall*, *F₁-score* and *accuracy* to evaluate the anomaly detection performance. The recall-oriented understudy for gisting evaluation (ROUGE) (Lin 2004) is a software package and metric set designed to assess the quality of generated text by comparing it with ground truth text. In our evaluation of DABL's ability to interpret the cause of anomalies, we utilize *ROUGE-2* and *ROUGE-L* metrics. We conduct each experiment five times and report the mean results.

Quantitative Results

Anomaly Detection To evaluate the model's **generalization ability**, we conduct experiments on the test dataset \mathcal{D}_1 . The results are shown in Table 1, with the best outcomes highlighted in bold. Our DBAL achieves the highest precision, F₁-score, and accuracy, with both the F₁-score and accuracy exceeding 90%. Although SENSE-BERT attains the best recall, it has the lowest precision and accuracy.

Dataset	ROUGE-2(%)			ROUGE-L(%)		
	Prec.	Rec.	F ₁	Prec.	Rec.	F ₁
\mathcal{D}_1	74.48	74.49	74.32	76.29	76.11	76.02
\mathcal{D}_2	84.92	84.61	84.54	86.96	86.66	86.56

Table 3: The results of DABL in interpreting the causes of anomalies.


Below is a business process trace, with each executed activity separated by a comma: (Add personal data, Procure work equipment, Select necessary work equipment, Set up access rights hardware and software). Is this trace normal or anomalous? 


GPT-3.5 Without further context or information, it is difficult to determine if the given business process trace is normal or anomalous.

GPT-4 Without knowing the standard sequence of activities or having additional context, it's difficult to say definitively whether this trace is normal or anomalous.

GLM-3 This sequence of activities seems logical and typically follows an onboarding process in many organizations. It would be considered normal if this is the standard procedure for equipping new employees with the necessary tools and access to perform their job.

Llama 2 A business process trace is a sequence of activities that are executed in a business process. Each executed activity is separated by a comma in the trace.

DABL  The trace is anomalous.

What's the origin of the anomaly? 



DABL  The activity 'Procure work equipment' is executed too early, it should be executed after 'Select necessary work equipment'.

Figure 5: Qualitative results of various LLMs on an *early* anomaly.

Compared to other methods, DBAL maintains a balanced precision and recall. SENSE-SVM exhibits limited sensitivity to anomalies, potentially overlooking many anomalies, thereby achieving high precision but markedly low recall. Conversely, SENSE-BERT demonstrates excessive sensitivity, resulting in numerous false alarms, thus yielding low precision but high recall. These results demonstrate that DBAL possesses a superior generalization ability for detecting anomalies in unseen processes.

We conduct experiments on dataset \mathcal{D}_2 to evaluate if the methods can **learn the given processes**. The results are presented in Table 2. Compared to the experiments on dataset \mathcal{D}_1 , the precision of each method increases. This improvement is due to the incorporation of normal process behavior in the training dataset, thereby reducing false alarms. Consequently, the F₁-score and accuracy also show significant improvements. Our DABL method stands out, achieving the best performance across all metrics, with the F₁-score and accuracy showing improvements of 21.79% and 16.75%, respectively, over the best existing semantic-based anomaly detection method (i.e., SENSE-BERT). These results demonstrate that DABL can learn the given processes effectively, allowing users to further fine-tune DABL on their own processes to obtain a customized model.

Interpretation of the Cause of Anomalies Table 3 shows the results of DABL in interpreting anomalies' causes. On


Here is a trace of a business process, with each performed activity separated by a comma: (Ship goods, Contact container port, Contact airlines, Prepare delivery note). Is this trace normal or anomalous? 


GPT-3.5 Based on the provided business process trace, it is difficult to determine if it is normal or anomalous without further context or information.

GPT-4 As an AI, I don't have the exact context or knowledge of the specific rules of the business process. However, based on the general order of operations in shipping and logistics, this trace seems normal.

GLM-3 Without additional context or a definition of what is considered normal in this specific process, we can only say that the given trace is not inherently anomalous based on the activities listed.

Llama 2 What is the difference between a trace and a log? A trace is a record of all the activities performed in a business process. A log is a record of all the activities performed in a business process.

DABL  The trace is anomalous.

What makes this trace anomalous? 


DABL  The activity 'Contact airlines' is mutually exclusive with the activity 'Contact container port', meaning they should not be executed within the same process instance.

Figure 6: Qualitative results of various LLMs on an *exclusion* anomaly.

dataset \mathcal{D}_1 , both ROUGE-2 and ROUGE-L scores are relatively high, indicating that DABL performs well in identifying anomalies' causes, even for processes not included in the training data. For dataset \mathcal{D}_2 , DABL exhibits better performance because the normal behaviors of the processes are well-represented in the training data. Furthermore, the slight difference between recall and precision suggests that the model maintains a good balance. These results demonstrate that DABL is effective at interpreting anomalies' causes regarding both bigrams and longest common subsequences.

However, these results may be underestimated because the cause of an anomaly can be interpreted in various ways. For example, for a desired trace $\langle A, B, C, D, E \rangle$, the anomaly $\langle A, B, E, C, D \rangle$ can be interpreted as "The activities 'C' and 'D' are executed too late, they should be executed after 'B'" and "The activity 'E' is executed too early, it should be executed after 'D'." Nevertheless, we only provide one reference answer to calculate ROUGE-2 and ROUGE-L scores.

Qualitative Examples

Fig. 5 and Fig. 6 illustrate the performance of existing LLMs and DABL on *early* and *exclusion* anomalies, respectively. Due to the absence of business process-related knowledge, existing LLMs may not offer definitive answers. In contrast, DABL is proficient in identifying anomalies and interpreting their causes, thanks to the integration of a vast array of process models from diverse domains in the training data.

Real-world Application

Travel Permit We apply our DABL on a real-world travel permit log from the BPI 2020 challenge (van Dongen 2020), which captures data on work trips conducted by university employees. The process flow involves the request for and approval of a travel permit, the trip itself, a subsequent travel declaration, as well as associated reimbursements.

ID	Example Trace	Output Cause
A1	(ST, PSE, PAA, PAB, PAS, PFAD, DSE, DAA, DAB, DFAS, RP, PH, ET)	The activity 'Start trip' should not be executed.
A2	(PSE, ST, PAA, PAB, PAS, PFAD, DSAE, ET)	The activity 'Start trip' is executed too early, it should be executed after 'Permit FINAL_APPROVED by DIRECTOR'.
A3	(PSE, PAS, ST, PFAD, PRM, ET)	The activity 'Permit REJECTED by MISSING' is mutually exclusive with the activity 'Start trip', meaning they should not be executed within the same process instance.
A4	(ST, ET, PSE, PFAS, DSE, DFAS, RP, PRM, PH)	The activity 'Permit REJECTED by MISSING' is mutually exclusive with the activities 'Payment Handled' and 'Declaration FINAL_APPROVED by SUPERVISOR', meaning they should not be executed within the same process instance.
A5	(PSE, PAA, RSE, RAA, RAB, RFAS, RP, PAB, PFAS, PH, ST, ET)	The activities 'Permit APPROVED by BUDGET OWNER' and 'Permit FINAL_APPROVED by SUPERVISOR' are executed too late, they should be executed before 'Request For Payment APPROVED by BUDGET OWNER'.

PSE: Permit SUBMITTED by EMPLOYEE; **PAA:** Permit APPROVED by ADMINISTRATION; **PAB:** Permit APPROVED by BUDGET OWNER; **PAS:** Permit APPROVED by SUPERVISOR; **PFAS:** Permit FINAL_APPROVED by SUPERVISOR; **PFAD:** Permit FINAL_APPROVED by DIRECTOR; **PRM:** Permit REJECTED by MISSING; **RSE:** Request For Payment SUBMITTED by EMPLOYEE; **RRA:** Request For Payment REJECTED by ADMINISTRATION; **RRE:** Request For Payment REJECTED by EMPLOYEE; **RAA:** Request For Payment APPROVED by ADMINISTRATION; **RAB:** Request For Payment APPROVED by BUDGET OWNER; **RFAS:** Request For Payment FINAL_APPROVED by SUPERVISOR; **DSAE:** Declaration SAVED by EMPLOYEE; **DSE:** Declaration SUBMITTED by EMPLOYEE; **DRA:** Declaration REJECTED by ADMINISTRATION; **DRE:** Declaration REJECTED by EMPLOYEE; **DAA:** Declaration APPROVED by ADMINISTRATION; **DAB:** Declaration APPROVED by BUDGET OWNER; **DFAS:** Declaration FINAL_APPROVED by SUPERVISOR; **RP:** Request Payment; **PH:** Payment Handled; **SR:** Send Reminder; **ST:** Start Trip; **ET:** End Trip.

Table 4: Irregularity patterns identified in the travel permit log.

ID	Example Trace	Output Cause
A1	(CF, SF, IFN, AP, P, P, P, P)	The activities 'Payment' and 'Payment' are reworked after 'Payment'.
A2	(CF, P, P, SF, IFN, AP, P)	The activities 'Payment' and 'Payment' are executed too early, they should be executed after 'Add penalty'.
A3	(CF, SF, IFN, RRAP, IDAP, SAP, AP, P)	The activity 'Receive Result Appeal from Prefecture' is executed too early, it should be executed after 'Send Appeal to Prefecture'.

CF: Create Fine; **SF:** Send Fine; **IFN:** Insert Fine Notification; **AP:** Add Penalty; **P:** Payment; **RRAP:** Receive Result Appeal from Prefecture; **IDAP:** Insert Date Appeal to Prefecture; **SAP:** Send Appeal to Prefecture.

Table 5: Irregularity patterns identified in the road traffic fine management log.

This log contains 7,065 traces with 1,478 variants (traces with different orders of activities are executed). DABL identifies 562 anomalous variants. These detected anomalies reveal some interesting irregularity patterns illustrated in Table 4. These irregularity patterns include: trips starting before a permit is properly handled (A1), approved (A2), or even rejected (A3); the declaration being finally approved by a supervisor and payment handled despite the permit being rejected (A4); and requests for payment being approved before the permit is approved (A5).

Road Traffic Fine Management We apply our DABL on another real-world event log from an information system managing road traffic fines (de Leoni and Mannhardt 2015), which captures the road traffic fine management process. The process flow involves the creation of a fine, appeal to the prefecture, addition of penalties, and fine payment.

This log contains 150,370 traces with 231 variants. DABL identifies 56 anomalous variants, which reveal some interesting irregularity patterns as detailed in Table 5. The examples illustrate irregularity patterns where the fine is repeatedly paid (A1), the fine is paid before the penalty is added (A2), and the result appeal from the prefecture is received before the appeal is sent to the prefecture (A3).

Limitations

Despite collecting a significant amount of data to fine-tune LLMs for detecting business process anomalies and achieving good performance, it is important to acknowledge some limitations of DABL.

First, the precision of the collected process model may not be equal to 1, which means it might allow traces not observed in the original process. This discrepancy could impact the quality of the training data, posing a limitation for DABL. Second, our open-source model may struggle with traces that, while irregular from a common-sense perspective, are normal within their customized processes. Finally, when multiple anomaly types occur within a trace, DABL can still classify the trace as an anomaly but can only interpret one type of anomaly at a time. After fixing the interpreted anomaly and re-inputting the trace into DABL, the model will identify and interpret another anomaly type.

Conclusion

In this paper, we introduce DABL, a novel semantic business process anomaly detection model leveraging LLMs. Trained on 143,137 real-world process models from various domains, DABL excels at zero-shot detection of semantic anomalies and interprets their causes in natural language. Extensive experiments demonstrate DABL’s generalization ability, allowing users to detect anomalies in their own datasets without additional training.

Acknowledgments

This work is supported by China National Science Foundation (Granted Number 62072301). This work is also partially supported by the Program of Technology Innovation of the Science and Technology Commission of Shanghai Municipality (Granted No.22DZ1100103).

References

- Aalst, W. v. d.; Buijs, J.; and Dongen, B. v. 2011. Towards improving the representational bias of process mining. In *International Symposium on Data-Driven Process Discovery and Analysis*, 39–54. Springer.
- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Berti, A.; Van Zelst, S. J.; and van der Aalst, W. 2019. Process mining for python (PM4Py): bridging the gap between process-and data science. *arXiv preprint arXiv:1905.06169*.
- Caspary, J.; Rebmann, A.; and van der Aa, H. 2023. Does this make sense? machine learning-based detection of semantic anomalies in business processes. In *International Conference on Business Process Management*, 163–179. Springer.
- Chinosi, M.; and Trombetta, A. 2012. BPMN: An introduction to the standard. *Computer Standards & Interfaces*, 34(1): 124–134.
- Chklovski, T.; and Pantel, P. 2004. Verbocean: Mining the web for fine-grained semantic verb relations. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, 33–40.
- De Koninck, P.; vanden Broucke, S.; and De Weerd, J. 2018. act2vec, trace2vec, log2vec, and model2vec: Representation Learning for Business Processes. In Weske, M.; Montali, M.; Weber, I.; and vom Brocke, J., eds., *Business Process Management*, 305–321. Cham: Springer International Publishing. ISBN 978-3-319-98648-7.
- de Leoni, M. M.; and Mannhardt, F. 2015. Road Traffic Fine Management Process.
- Dettmers, T.; Pagnoni, A.; Holtzman, A.; and Zettlemoyer, L. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Dumas, M.; Rosa, L. M.; Mendling, J.; and Reijers, A. H. 2018. *Fundamentals of business process management*. Springer.
- Ebrahim, M.; and Golpayegani, S. A. H. 2022. Anomaly detection in business processes logs using social network analysis. *J. Comput. Virol. Hacking Tech.*, 1–13.
- Egersdoerfer, C.; Zhang, D.; and Dai, D. 2023. Early exploration of using chatgpt for log-based anomaly detection on parallel file systems logs. In *Proceedings of the 32nd International Symposium on High-Performance Parallel and Distributed Computing*, 315–316.
- Elhafsi, A.; Sinha, R.; Agia, C.; Schmerling, E.; Nesnas, I. A.; and Pavone, M. 2023. Semantic anomaly detection with large language models. *Autonomous Robots*, 47(8): 1035–1055.
- Gu, Z.; Zhu, B.; Zhu, G.; Chen, Y.; Tang, M.; and Wang, J. 2024. Anomalygpt: Detecting industrial anomalies using large vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 1932–1940.
- Guan, W.; Cao, J.; Gu, Y.; and Qian, S. 2023. GRASPED: A GRU-AE Network Based Multi-Perspective Business Process Anomaly Detection Model. *IEEE Transactions on Services Computing*, 1–13.
- Guan, W.; Cao, J.; Gu, Y.; and Qian, S. 2024a. GAMA: A multi-graph-based anomaly detection framework for business processes via graph neural networks. *Information Systems*, 102405.
- Guan, W.; Cao, J.; Qian, S.; and Gao, J. 2024b. LogLLM: Log-based Anomaly Detection Using Large Language Models. *arXiv preprint arXiv:2411.08561*.
- Guan, W.; Cao, J.; Yao, Y.; Gu, Y.; and Qian, S. 2024c. COMB: Interconnected Transformers-Based Autoencoder for Multi-Perspective Business Process Anomaly Detection. In *2024 IEEE International Conference on Web Services (ICWS)*, 1115–1124. IEEE.
- Guan, W.; Cao, J.; Zhao, H.; Gu, Y.; and Qian, S. 2024d. Survey and Benchmark of Anomaly Detection in Business Processes. *IEEE Transactions on Knowledge and Data Engineering*, 1–23.
- Hoaglin, D. C.; and Welsch, R. E. 1978. The hat matrix in regression and ANOVA. *The American Statistician*, 32(1): 17–22.
- Huo, S.; Völzer, H.; Reddy, P.; Agarwal, P.; Isahagian, V.; and Muthusamy, V. 2021. Graph Autoencoders for Business Process Anomaly Detection. In *Business Process Management - 19th International Conference, BPM*, volume 12875 of *Lecture Notes in Computer Science*, 417–433. Springer.
- Junior, S. B.; Ceravolo, P.; Damiani, E.; Omori, N. J.; and Tavares, G. M. 2020. Anomaly Detection on Event Logs with a Scarcity of Labels. In *2020 2nd International Conference on Process Mining (ICPM)*, 161–168.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Ko, J.; and Comuzzi, M. 2021. Detecting anomalies in business process event logs using statistical leverage. *Information Sciences*, 549: 53 – 67.
- Ko, J.; and Comuzzi, M. 2022. Keeping our rivers clean: Information-theoretic online anomaly detection for streaming business process events. *Information Systems*, 104.
- Krajsic, P.; and Franczyk, B. 2021. Semi-supervised anomaly detection in business process event data using self-attention based classification. *Procedia Computer Science*, 192: 39–48.
- Leemans, S. J. J.; Fahland, D.; and van der Aalst, W. M. P. 2018. Scalable process discovery and conformance checking. *Softw. Syst. Model.*, 17(2): 599–631.
- Li, Y.; Wang, H.; Yuan, S.; Liu, M.; Zhao, D.; Guo, Y.; Xu, C.; Shi, G.; and Zuo, W. 2023. Myriad: Large multimodal model by applying vision experts for industrial anomaly detection. *arXiv preprint arXiv:2310.19070*.
- Lin, C.-Y. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, 74–81.

- Liu, J.; Zhang, C.; Qian, J.; Ma, M.; Qin, S.; Bansal, C.; Lin, Q.; Rajmohan, S.; and Zhang, D. 2024. Large Language Models can Deliver Accurate and Interpretable Time Series Anomaly Detection. *arXiv preprint arXiv:2405.15370*.
- Lu, K.; Fang, X.; and Fang, N. 2022. PN-BBN: A Petri Net-Based Bayesian Network for Anomalous Behavior Detection. *Mathematics*, 10(20).
- Mikolov, T.; Chen, K.; Corrado, G.; and Dean, J. 2013. Efficient Estimation of Word Representations in Vector Space. In Bengio, Y.; and LeCun, Y., eds., *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*.
- Niro, A.; and Werner, M. 2024. Detecting Anomalous Events in Object-centric Business Processes via Graph Neural Networks. *arXiv preprint arXiv:2403.00775*.
- Nolle, T.; Luetzgen, S.; Seeliger, A.; and Mühlhäuser, M. 2022. BINet: Multi-perspective business process anomaly classification. *Inf. Syst.*, 103: 101458.
- Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744.
- Pauwels, S.; and Calders, T. 2019a. An Anomaly Detection Technique for Business Processes Based on Extended Dynamic Bayesian Networks. In *Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing, SAC '19*, 494–501. New York, NY, USA: Association for Computing Machinery. ISBN 9781450359337.
- Pauwels, S.; and Calders, T. 2019b. Detecting Anomalies in Hybrid Business Process Logs. *SIGAPP Appl. Comput. Rev.*, 19(2): 18–30.
- Qi, J.; Huang, S.; Luan, Z.; Yang, S.; Fung, C.; Yang, H.; Qian, D.; Shang, J.; Xiao, Z.; and Wu, Z. 2023. Log-gpt: Exploring chatgpt for log-based anomaly detection. In *2023 IEEE International Conference on High Performance Computing & Communications, Data Science & Systems, Smart City & Dependability in Sensor, Cloud & Big Data Systems & Application (HPCC/DSS/SmartCity/DependSys)*, 273–280. IEEE.
- Rullo, A.; Guzzo, A.; Serra, E.; and Tirrito, E. 2020. A framework for the multi-modal analysis of novel behavior in business processes. In *Intelligent Data Engineering and Automated Learning—IDEAL 2020: 21st International Conference, Guimarães, Portugal, November 4–6, 2020, Proceedings, Part I 21*, 51–63. Springer.
- Sarda, K.; Namrud, Z.; Rouf, R.; Ahuja, H.; Rasoloveicy, M.; Litoiu, M.; Shwartz, L.; and Watts, I. 2023. Adarma auto-detection and auto-remediation of microservice anomalies by leveraging large language models. In *Proceedings of the 33rd Annual International Conference on Computer Science and Software Engineering*, 200–205.
- Sarno, R.; Sinaga, F.; and Sungkono, K. R. 2020. Anomaly detection in business processes using process mining and fuzzy association rule learning. *J. Big Data*, 7: 1–19.
- Sinaga, F.; and Sarno, R. 2016. Business process anomaly detection using multi-level class association rule learning. *IPTEK Journal of Proceedings Series*, 2(1).
- Sola, D.; Warmuth, C.; Schäfer, B.; Badakhshan, P.; Rehse, J.-R.; and Kampik, T. 2022. SAP Signavio Academic Models: a large process model dataset. In *International Conference on Process Mining*, 453–465. Springer.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- van der Aa, H.; Rebmann, A.; and Leopold, H. 2021. Natural language-based detection of semantic execution anomalies in event logs. *Information Systems*, 102: 101824.
- Van Der Aalst, W. M.; Van Hee, K. M.; Ter Hofstede, A. H.; Sidorova, N.; Verbeek, H.; Voorhoeve, M.; and Wynn, M. T. 2011. Soundness of workflow nets: classification, decidability, and analysis. *Formal aspects of computing*, 23: 333–363.
- van Dongen, B. 2020. BPI Challenge 2020.
- van Zelst, S. J.; and Leemans, S. J. 2020. Translating Workflow Nets to Process Trees: An Algorithmic Approach. *Algorithms*, 13(11): 279.
- Vertuam Neto, R.; Tavares, G.; Ceravolo, P.; and Barbon, S. 2021. On the Use of Online Clustering for Anomaly Detection in Trace Streams. In *XVII Brazilian Symposium on Information Systems, SBSI 2021*. New York, NY, USA: Association for Computing Machinery. ISBN 9781450384919.
- Weske, M.; Decker, G.; Dumas, M.; La Rosa, M.; Mendling, J.; and Reijers, H. 2020. Model collection of the business process management academic initiative. *DOI: https://doi.org/10.5281/zenodo.3758705*.
- Wu, J.; and Hooi, B. 2023. Fake News in Sheep’s Clothing: Robust Fake News Detection Against LLM-Empowered Style Attacks. *arXiv preprint arXiv:2310.10830*.
- Yang, L.; McClean, S.; Donnelly, M.; Khan, K.; and Burke, K. 2020. Analysing Business Process Anomalies Using Discrete-time Markov chains. In *2020 IEEE 22nd International Conference on High Performance Computing and Communications; IEEE 18th International Conference on Smart City; IEEE 6th International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*, 1258–1265.
- Zeng, A.; Liu, X.; Du, Z.; Wang, Z.; Lai, H.; Ding, M.; Yang, Z.; Xu, Y.; Zheng, W.; Xia, X.; et al. 2022. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*.