

SHAining on Process Mining: Explaining Event Log Characteristics Impact on Algorithms

Andrea Maldonado^{*†‡}, Christian M. M. Frey^{||}, Sai Anirudh Aryasomayajula^{*},
Ludwig Zellner^{*}, Stephan A. Fahrenkrog-Petersen[§], Thomas Seidl^{*†}

^{*}Database Systems and Data Mining, Ludwig Maximilian University of Munich, Germany

[†]Munich Center for Machine Learning, Germany

[‡]School of Engineering and Design, Technical University of Munich, Germany

^{||}Machine Learning Lab, University of Technology Nuremberg, Germany

[§]University of Liechtenstein, Liechtenstein

andrea.maldonado@tum.de, christian.frey@utn.de, anirudhsai027@gmail.com,
zellner@dbs.ifi.lmu.de, stephan.fahrenkrog@uni.li, seidl@dbs.ifi.lmu.de

Abstract—Process mining aims to extract and analyze insights from event logs, yet algorithm metric results vary widely depending on structural event log characteristics. Existing work often evaluates algorithms on a fixed set of real-world event logs but lacks a systematic analysis of how event log characteristics impact algorithms individually. Moreover, since event logs are generated from processes, where characteristics co-occur, we focus on associational rather than causal effects to assess how strong the overlapping individual characteristic affects evaluation metrics without assuming isolated causal effects, a factor often neglected by prior work. We introduce SHAining, the first approach to quantify the marginal contribution of varying event log characteristics to process mining algorithms’ metrics. Using process discovery as a downstream task, we analyze over 22,000 event logs covering a wide span of characteristics to uncover which affect algorithms across metrics (e.g., fitness, precision, complexity) the most. Furthermore, we offer novel insights about how the value of event log characteristics correlates with their contributed impact, assessing the algorithm’s robustness.

Index Terms—Explainability, Shapley Value, Feature Contribution, Algorithm Evaluation, Process Discovery

I. INTRODUCTION

Process mining (PM) techniques are widely used to extract actionable insights from event logs across various domains such as healthcare, manufacturing, and finance. Among these techniques, process discovery, which constructs process models from event logs, has received substantial attention [1], [2]. However, evaluating the quality of discovered models remains a challenge, particularly due to the heterogeneity in event log characteristics [3], [4], [5]. Event logs can differ significantly in structural properties. For instance, healthcare processes typically yield highly variable traces where most cases exhibit similar behavior [6], whereas logs from structured domains like manufacturing often show recurring patterns with a few dominant variants.

Yet, not all process discovery algorithms are equally equipped to handle variability in event log characteristics. Previous studies have related algorithmic metrics to either structural descriptions of process models [7], [8] or event log characteristics [9] via statistical approaches such as linear

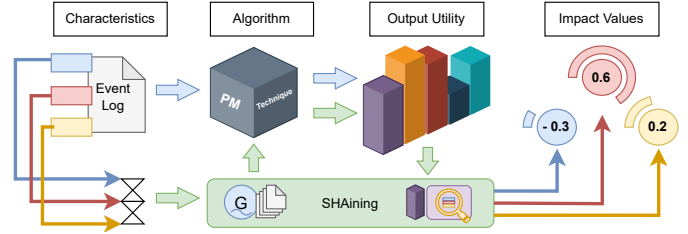


Fig. 1: Studying the impact of log characteristics on algorithm output with SHAining.

regression, assuming a fixed dataset and a fixed model. Moreover, these approaches assume properties like homoscedasticity and feature independence, which are rarely satisfied in practice. Prior work has examined how variation in event log feature values influences overall algorithm evaluation metrics; the question of “why” a given technique performs well or poorly on a specific log remains largely unaddressed.

In this paper, we present the first systematic and explainable analysis of how overlapping event log characteristics individually contribute to impact on process mining algorithms. Rather than optimizing algorithm selection for a downstream task, our goal is to explain metric variation by treating process mining techniques as black-boxes and analyzing how structural log characteristics affect evaluation outcomes. To this end, we introduce SHAining, depicted as the green bottom box in Figure 1, which quantifies the impact of characteristics, highlighted in {blue, red, yellow}, to algorithm metrics in a technique-agnostic way. We generate over 22,000 intentional logs with varied structural properties using GEDI [3], and evaluate multiple process mining algorithms with standard metrics such as fitness, precision, and F-score. We apply a Shapley value analysis [10], [11] originating from game theory that fairly attributes the impact across all possible feature value combinations. We identify which feature values, e.g. entropy or variant diversity, affect an algorithm’s performance the most. In our analysis, we study the impact of event log characteristics on algorithm metrics using process discovery as a downstream task. We demonstrate that contributions of

event log characteristics to algorithm metrics vary by the type of structural characteristic and by their specific values. For instance, high activity entropy degrades the output of some algorithms, revealing key trade-offs in algorithm suitability.

Overall, our key contributions are: (1) Extends Shapley analysis to distributions defined by generative models conditioned on meta-feature configurations. (2) Enables explainability of how distributional properties influence downstream model behavior. (3) Large-scale analysis on 22,000+ synthetic event logs, revealing consistent feature-metric interactions.

II. RELATED WORK

Understanding the impact of event log characteristics on process discovery (PD) algorithms has become a growing area of interest in PM. A prime example is the Process Discovery Contest [12] which annually compares PD algorithms to determine the best-performing approach.

Explainability. Shapley value analysis is applied across various domains. Stevens et al. [13] apply Shapley value analysis to compare interpretable models and post-hoc explainability methods for predicting loan application outcomes. Similarly, Pishgar et al. [14] use Shapley values to quantify event attribute importance in predicting COVID-19 mortality. Heskes et al. [15] extend Shapley-based explanations by incorporating causal knowledge, highlighting the importance of distinguishing between association and causation when interpreting model outputs. Although these works demonstrate the usefulness of Shapley values for explainability in PM scenarios, they focus on predictive tasks rather than directly evaluating PD algorithms or explaining their behavior on specific log structures. A closer line of research is the ProReco framework [7], a recommender system that predicts performance metrics for different PD algorithms based on log characteristics and user preferences. While ProReco uses SHAP (SHapley Additive exPlanations) [16] to provide explanations, it applies SHAP on a surrogate model rather than the PD algorithms themselves. Additionally, it relies on heuristic imputation for incomplete data, and is limited to a much smaller synthetic dataset.

A recent critique [17] argues that Shapley values can invalidate intuitions and induce false trust by assigning high scores to irrelevant features. They propose applying formal methods, such as abductive/contrastive explanations (AXps/CXps) [18]. However, this critique targets the usage of Shapley analysis on a fixed per-instance level, i.e., explaining individual predictions. DShapley [19] extends the classical Shapley value to quantify the contribution of features w.r.t. the expected utility over an underlying data distribution, rather than a fixed dataset. In our work, we compute Shapley values on meta-features that govern data generation to analyze impacts on distributions (cf. [19]) of algorithm performance, not on individual predictions of a *fixed* model on a *fixed* input. Hence, SHAining differs fundamentally from the fixed data scenario that AXps/CXps address. Therefore, Shapley values over distributions [19] arise as the most suitable for our analysis of associational trends across generative configurations.

Benchmark studies. Several benchmark studies compare PD algorithms across diverse logs and quality metrics. Augusto et al. [20] evaluate multiple algorithms across 24 real-world logs using nine different evaluation criteria, providing insights into the strengths and weaknesses of these algorithms. Van den Broucke et al. [9] analyze how structural log characteristics influence PD outcomes using metrics such as fitness and precision. Yet, their reliance on linear regression introduces strong assumptions such as feature independence and homoscedasticity that are often breached in practice due to the complex dependencies and variance inherent in real-life event logs. A complementary direction is taken by Andree et al. [21], who compare PD algorithms based on their ability to reproduce control-flow patterns. However, their focus is on what patterns are captured, rather than how structural log characteristics influence evaluation outcomes, a gap our work addresses directly.

Generation. Recent work for data generation has tackled the scarcity and imbalance of benchmark logs. Maldonado et al. [3] introduce GEDI, a framework for generating logs with diverse structural characteristics, enabling controlled experimentation. Similarly, Jouck et al. [5] propose a generator based on random sampling from predefined process populations. Janssenwillen et al. [22] explore how noise injection affects fitness and precision in PD evaluations, particularly for bias estimation. While these efforts improve empirical validity, they stop short of explicitly analyzing how individual log characteristics affect evaluation outcomes.

III. PRELIMINARIES

Event Log. An event $e \in \mathcal{E}$ represents a step in a process and is characterized by $e := (c, a, t)$, with a case identifier $c \in \mathcal{C}$, an activity $a \in \mathcal{A}$ describing the type of event, and a timestamp $t \in \mathcal{T}$ when the event occurred. A sequence of events is called a trace $\sigma := \langle e_1, e_2, \dots, e_n \rangle$ and groups events that belong to the same case within the process. Within a trace, events are ordered based on their timestamp. Finally, we define a multiset of traces as an event log L . We denote the universe of all possible event logs as \mathcal{L} .

Event Log characteristics. Event logs capture the execution of processes, which can exhibit a wide spectrum of characteristics, ranging from highly structured to highly variable and complex workflows. In our work, we quantify these characteristics using *event log features* $\mathcal{F} := \{F_1, \dots, F_n\}$, which describe different aspects related to traces, variants, activities, and events, such as trace lengths, and activity frequencies. Following Maldonado et al. [23], feature extraction is formalized as a function $f_e : \mathcal{L} \rightarrow \mathbb{R}^n$, mapping an event log $L \in \mathcal{L}$ to an n -dimensional feature vector $f := (f_1, \dots, f_n)$, where each $f_i \in \mathbb{R}$ represents a specific feature value.

Shapley Value Analysis. Shapley value analysis [10] originates from game theory, where we define a set of players $N = \{1, \dots, n\}$ and a characteristic function $v : 2^N \rightarrow \mathbb{R}$ that assigns to each coalition of players $S \subseteq N$ a real number

$v(S)$ ¹. This concept has been widely adapted in machine learning to interpret model predictions, where "players" correspond to input features, and $v(S)$ represents the model's output when only the subset S of features is considered. The Shapley value thus offers a rigorous, fair method to evaluate the influence of each feature on a model's prediction. In our work, we examine the impacts of meta-features on the data-generating process, rather than on fixed input datasets.

IV. SHAINING - SHAPLEY VALUES FOR RELATING PROCESS MINING TASKS AND EVENT LOG FEATURES

We apply Shapley values [10], [11], [16] for a higher-order problem across the space of data-generating processes. Concretely, our players represent features in the data generation process itself. For a well-defined characteristic function, we apply a Shapley analysis, where the *players* are (meta-)features that generate logs, and the *game* is defined by a model's performance. More concretely, we answer the question: *What is the marginal value of including feature F_i in the log generation process, in terms of the downstream utility of a PM technique?* Therefore, we are not comparing feature subsets for a *fixed* model on a *fixed* dataset, but rather assessing the marginal effect of a feature on the generative process, and through it, on a model's evaluation results.

Our approach consists of four steps: **(i) Feature Combination.** **(ii) Event log Generation.** **(iii) Process Mining.** **(iv) Shapley Value Computation.** The first three steps are repeatedly applied for a robust Shapley Value analysis across a diverse set of event logs being processed by a black-box model. An overview of SHAINing is summarized in Figure 2. In this section, we elaborate on each step individually:

(i) Feature Combination. As in section III, let the domain of (meta-)features be \mathcal{F} . We define an n -dimensional configuration vector as a specific realization of features: A) First, we initialize 1-dimensional vectors, each as $f_{S_i} := \{(F_i = f_{S_i}) | F_i \in \mathcal{F}\}$, each containing a value for the i -th feature; B) Later, any $(|S_i| + 1)$ -dimensional vector is $f_{S_i \oplus S_j} := f_{S_i} \oplus f_{S_j}$ composed of two $|S_i|$ -dimensional feature vectors $f_{S_i}, f_{S_j} \in \mathbb{R}^{|S|}$, of a subset $S \subseteq \mathcal{F}$ of features, where \oplus is the subspace join operator, resulting in the smallest common superspace. In our evaluation, we utilize a correlation analysis on features to define \mathcal{F} (see Section V).

(ii) Event log Generation. We leverage an event log generator to produce different event logs targeting respective feature configurations. Generally, assuming that we have access to a prior over all event logs in the universe, denoted as $L \sim p(L)$, from which we can sample infinitely many logs. In that case, we can analyze an infinite number of feature combinations and their resulting metric distributions. In our work, we assume $p(L|f_S)$ to be a parameterized prior distribution over event logs with a feature configuration vector f_S controlling the feature distributions. The log generator \mathcal{G} is defined as a sampling function $\mathcal{G} : f_S \mapsto L \sim p(L|f_S)$. As a generator, we

leverage GEDI [3] as prior, resulting in event logs that align features to target feature values in the configuration vector.

(iii) Process Mining. We leverage a utility function $U(A, L_S)$ to measure a black-box model's quality in terms of evaluation metrics, when applying a fixed process mining algorithm A on the data L_S . We define:

$$v(S) := U(A, L_S) \quad (1)$$

In our evaluation, we apply Equation (1) for \mathcal{A} different black box models and \mathcal{M} metrics on each generated event log. Whenever a feature configuration could be generated, and the log passed the process mining analysis (steps i-iii), we evaluate whether the dimension of the configuration feature vector equals the number of players. If it does not, we (i) combine the current feature vectors pair-wise, otherwise we continue to (iv) Shapley Value Computation.

(iv) Shapley Value Computation. We define the functional Shapley value for a feature $f_i \in \mathcal{F}$ as [11], using the PM evaluation measurements from utility function $v(S)$. By averaging the marginal contributions across all coalitions, we get the final Shapley value per feature that quantifies its average impact on an algorithm's evaluation results:

$$\phi_i = \sum_{S \subseteq \mathcal{F} \setminus \{f_i\}} \frac{|S|!(|\mathcal{F}| - |S| - 1)!}{|\mathcal{F}|} [v(S \cup \{f_i\}) - v(S)], \quad (2)$$

Discussion. The following properties are fulfilled[11]:

Axiom 1 (Efficiency): The total gain of the feature set \mathcal{M} is distributed: $\sum_{i=0}^{|\mathcal{F}|} \phi_i = v(\mathcal{F})$

Axiom 2 (Symmetry): If two features $f_i, f_j \in \mathcal{F}$ contribute equally to all possible coalitions, i.e., $v(S \cup \{f_i\}) = v(S \cup \{f_j\})$, then their Shapley values are identical: $\phi_i = \phi_j$

Axiom 3 (Additivity): If two coalitions of features defined by gain functions v and w are combined, then the distributed gains correspond to gains derived from v and w , i.e., $\phi_i(v + w) = \phi_i(v) + \phi_i(w)$. Additionally, for any scalar $a \in \mathbb{R}$, we have $\phi_i(a \cdot v) = a\phi_i(v)$.

Axiom 4 (Null player): If it holds that $v(S \cup f_i) = v(S), \forall S \subseteq \mathcal{F} \setminus \{f_i\}$, then the feature m_i is called a null player with $\phi_i = 0$.

Axioms **Additivity** 3 and **Null Player** 4 are trivially satisfied: the former holds as a null player's marginal contributions to any coalition are zero; the latter directly follows from linearity in the second term. Our application satisfies the **Symmetry** 2 axiom as our utility function (cf. Equation (1)) does not distinguish between features that induce identical shifts in data distributions and lead to identical model performance. In such cases, their marginal contributions are the same across all permutations. Finally, despite our utility function $v(S)$ being based on a fixed model applied to non-fixed data sampled from a conditioned generative process, the **Efficiency** 1 axiom still holds due to the combinatorial structure used to compute each feature's marginal contribution across all permutations, as shown by rewriting Equation (2) as:

¹The condition $v(\emptyset) = 0$ holds per definition of Shapley values.

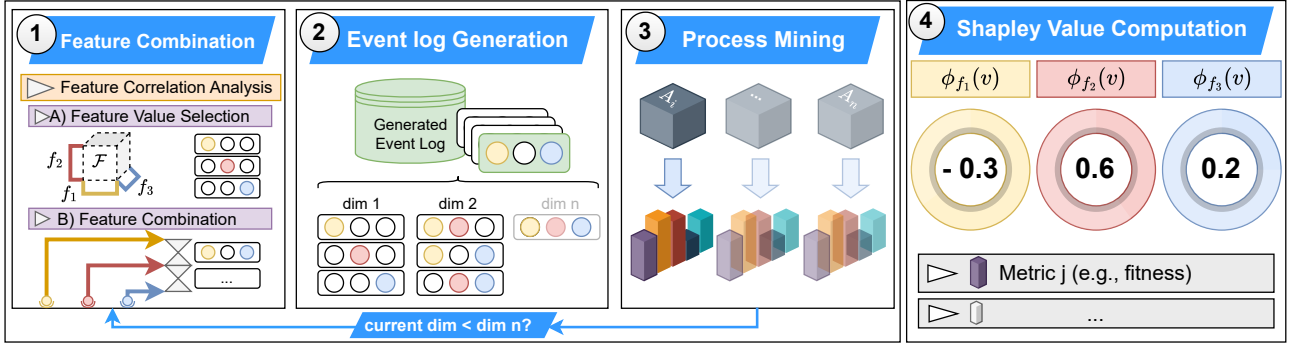


Fig. 2: The workflow of SHaining.

$$\begin{aligned}
 \sum_{i=0}^{|\mathcal{F}|} \phi_i &= \frac{1}{|\mathcal{F}|!} \sum_{\pi \in \text{Perm}(\mathcal{F})} \sum_{j=1}^{|\mathcal{F}|} [v(P_\pi(f_j) \cup \{f_j\}) - v(P_\pi(f_j))] \\
 &= \frac{1}{|\mathcal{F}|!} \sum_{\pi \in \text{Perm}(\mathcal{F})} [v(\mathcal{F}) - v(\emptyset)] = v(\mathcal{F}),
 \end{aligned}$$

where P_π denotes the set of players preceding the addition of feature f_i . When aggregating these contributions, the marginal effects follow a telescoping sum across all permutations spanning the entire value of $v(\mathcal{F})$. In our definition, even though datasets are different for each coalition, the utility function is a well-defined set function over $2^{\mathcal{F}} \rightarrow \mathbb{R}$ yielding a scalar value for all coalitions. Thus, even though $v(S)$ is induced by a different data distribution, the Shapley value ensures that efficiency is preserved in our setting.

V. SHAINING ON PROCESS DISCOVERY

We answer the following research questions:

- RQ1** Can we reveal which EL characteristics contribute the most to the algorithm's evaluation results?
- RQ2** Do EL characteristics' impacts on algorithm evaluation results change depending on feature values?
- RQ3** How can insights about the impact of EL characteristics support process miners on a specific downstream task?

Setup and Implementation Details. We run our framework on a Intel(R) Xeon(R) Platinum 8160 CPU @ 2.10GHz using 239Gi RAM. Our code is publicly available². Furthermore, each algorithm was tested under resource constraints with a maximum of 5 minutes and 19GB of disk storage per log.

Features. For our experiments, we selected a subset from FEEED [23] using a greedy approach. This method minimizes inter-feature correlation while maintaining representative coverage. As in Huang et al. [7], we use the Pearson correlation coefficient to identify representative meta-features from different groups, e.g., activity-/time-/trace-based features. We iteratively chose the meta-feature with the lowest average correlation and then added the feature with the smallest maximum correlation to the selected set. Using the elbow method, we set a cutoff point of $n=8$ meta-features. More

details on this pre-selection step are available in our repository. Moreover, we select the features with ten equidistant value samples. Thus, the feature set in our evaluation is $F = \{\text{aq1}, \text{nusa}, \text{saq1}, \text{ekbr3}, \text{rt5v}, \text{svo}, \text{tlkh}, \text{tlv}\}$, with:

- activities_q1 (aq1):** Lowest 25% (quartile) of activity counts in the log. Range: [1.0, 79.92].
- n_unique_start_activities (nusa):** Counts unique start activities in the log, indicating process diversity at the start. Range: [1.0, 6.56].
- start_activities_q1 (saq1):** Lowest 25% of start activity counts in the log. Range: [1.0, 174.79].
- evententropy_k_block_ratio_3 (ekbr3):** Normalized ratio of the 3-subsequence entropy in a log [24]. Range: [0.0, 4.37].
- ratio_top_5_variants (rt5v):** The proportion of traces in the top 5% most frequent variants. Range: [0.0, 0.38].
- skewness_variant_occurrence (svo):** Measures how unevenly process variants occur, showing if the distribution is balanced or not. Range: [1.54, 11.61].
- trace_len_kurtosis (tlkh):** Measures how much the trace lengths vary, indicating the concentration of items at the center. Range: [-0.97, 7.92].
- trace_len_variance (tlv):** Measures how much the lengths of different process variants vary. Range: [0.0, 138.7].

Datasets. Logs are generated during SHaining's generation step (see Section IV), with the number of logs resulting to $\sum_{k=1}^{k_{max}} \binom{n}{k} v^k$ where n is the number of features, v is the number of values per feature, and $k_{max} \leq n$ is the maximal number of features in one single configuration vector (cf. Section IV). In our evaluation, we set $k_{max}=3$ to three features. With $v=10$ values and $n=8$ features, the number of possible feature combinations yields 58,880 logs. We refer to our repository for a summary of the event log statistics.

PM Downstream Task. SHaining's modular architecture, decoupled from black-box mechanics, enables seamless extension to various PM tasks. To quantify Shapley Value insights in PM, we use *process discovery (PD)* as a representative downstream task. PD is a fundamental PM task that automatically generates process models from event logs. Our evaluation compares three well-established *PD algorithms*, selected from different discovery paradigms, as top-down vs. bottom-up, for

²<https://github.com/andreamalhera/SHaining/tree/icpm25>

thorough assessment. We use default hyperparameters from their original works. Inductive Miner (*IND*) [25] uses a top-down approach, recursively partitioning event logs to construct models. Integer Linear Programming Miner (*ILP*) [26] uses optimization to incrementally identify patterns bottom-up. Split Miner (*SPM*) [27] builds sound process models by examining directly-follows graphs and loops first.

For fair and robust analysis, we evaluate PD outputs using multiple *metrics*, as in [3] and [28]:

Quality metrics: *Fitness* (*Ft*), *precision* (*Pr*), *F-score* (*Fs*) [29]; higher values indicate better model-log alignment [30], [31].

Complexity metrics: *Size* (*Sz*), i.e. number of BPMN nodes, and *control-flow complexity* (*Cf*), i.e. degree of split gateway branching [32].

Performance metrics: *Execution time* (*Et*) and model soundness, i.e., behavioral correctness [33], [20].

All metric evaluation results and Shapley Values per feature-value combination are in our repository. Our remaining analysis focuses on the resulting Shapley Values. For example, a Shapley value of 3.4 for trace length kurtosis (*tlkh*) concerning *Sz* and *ILP* indicates this feature increases model size by 3.4 on average. Normalized, a 0.22 value means *tlkh* accounts for 22% of the total feature impact on this metric.

A. RQ1: Can we reveal which EL characteristics contribute the most to the algorithm's evaluation results?

We evaluate the empirical average marginal effect over multiple datasets, i.e., $\bar{\phi} = \frac{1}{k} \sum_{i=1}^k \phi_i$, to determine the relative impact of features across various utility functions and reveal the most influential ones. We use the *autorank* package [34], to rank the (meta-)features by their statistical significance. Figures 3 and 4 use critical difference diagrams to show feature contribution. Features connected by black bars show no statistically significant difference in rank.

Figure 3 shows the results for six metrics across all evaluated PD algorithms. Notably, metrics *Ft*, *Cf*, and *Sz* rank *aq1* as feature with the highest contribution as a low number of activities directly simplifies the process model, improving its alignment with the log while reducing complexity and size. Similarly, *Pr* and *Et* rank *svo* as the most impactful feature as a highly skewed distribution of process variants allows for the creation of simpler, more precise models that are faster to generate. As the harmonic mean of *Ft* and *Pr*, *Fs* expectedly ranks *aq1* and *svo* as statistically equally important. On the other end, features ranked last reflect differences between metric types. Quality metrics (*Ft*, *Pr*, *Fs*) and the performance metric *Et* rank *rt5v* are the least impactful as these metrics are more influenced by the overall distribution of variants rather than the proportion of rare variants. Although *rt5v* is similarly ranked for complexity metrics *Cf*, *Sz*, the lowest contribution is average by *ekbr3*, as a statistically unpredictable log does not always require a large or intricate model to represent it.

Figure 4 shows the perspective per PD algorithms across all evaluated metrics. The features *svo* and *aq1* show the highest contribution for all miners. For *IND* and *ILP*, the rankings

of *svo* and *aq1* show no significant difference, while *SPM* ranks *aq1* first and *svo* second. Likewise, *rt5v* has the lowest contribution across all evaluated metrics and algorithms. The difference for *SPM* likely stems from its strategy of explicitly examining directly-follows graphs and loops, a process more influenced by the number of activities (*aq1*) than variant distribution (*svo*). Conversely, *rt5v* has the lowest contribution because it accounts for a small portion of the log's behavior. More holistic features, as *svo* and *aq1*, are therefore on average more influential on results across all algorithms.

B. RQ2: Do EL characteristics' impacts on algorithm evaluation results change depending on feature values?

To assess how log feature values influence process discovery results, we discuss three aspects: the *robustness* of algorithms to feature value changes, the *correlation* between feature values and Shapley values, and the *feasibility* of producing results with a given algorithm under varying feature conditions.

Algorithm Robustness. To evaluate algorithm sensitivity to event log characteristics, Figure 5a plots the normalized mean and variance of Shapley values across features. The mean Shapley value (as in *RQ1*) assesses overall feature influence on each algorithm, while variance measures stability across feature combinations. Higher mean values indicate stronger, more consistent feature influence; higher variance signals uneven responses to different feature value configurations. All miners are affected by feature variations, but their response stability differs. *IND* shows the highest average Shapley Value for *Fs*, but exhibits low variance, confirming its robustness due to its structured, noise-tolerant design. In contrast, the *ILP* is 40% more sensitive to feature variations, despite having 7% lower average Shapley Values. This reflects *ILP* sensitivity to infeasible or rare behavior to create precise and generalizable models. *SPM* falls close to *IND*, balancing responsiveness and stability slightly better, likely due to its embedded filtering.

Correlations. To assess how log feature values influence the explanatory power of features on PD outcomes, we analyze Spearman correlations between feature values and their corresponding Shapley values across algorithms and metrics. Spearman does not assume linearity and captures monotonic relationships. Figure 5b shows statistically significant correlations ($p \leq 0.05$) between feature values and their Shapley values, indicating whether increasing a feature consistently amplifies or reduces its impact on an evaluation metric. For *ILP*, features like the first quartile of starting activities (*saq1*) show positive correlations with quality metrics (*Pr*, *Fs*) but negative correlations with complexity metrics (*Cf*, *Sz*). A regular process start (high *saq1*) leads to better quality and simpler models, as the algorithm can identify a few dominant entry points for the complex optimization problem. In contrast, higher values of entropy for activity triplets (*ekbr3*), dominance of most common variants (*rt5v*), quantity of start activities (*nusa*), variant occurrence skewness (*svo*), and trace length variance (*tlv*) negatively impact (*Pr*, *Fs*), while increasing (*Cf*, *Sz*) measurements. This leads to more complicated and less precise

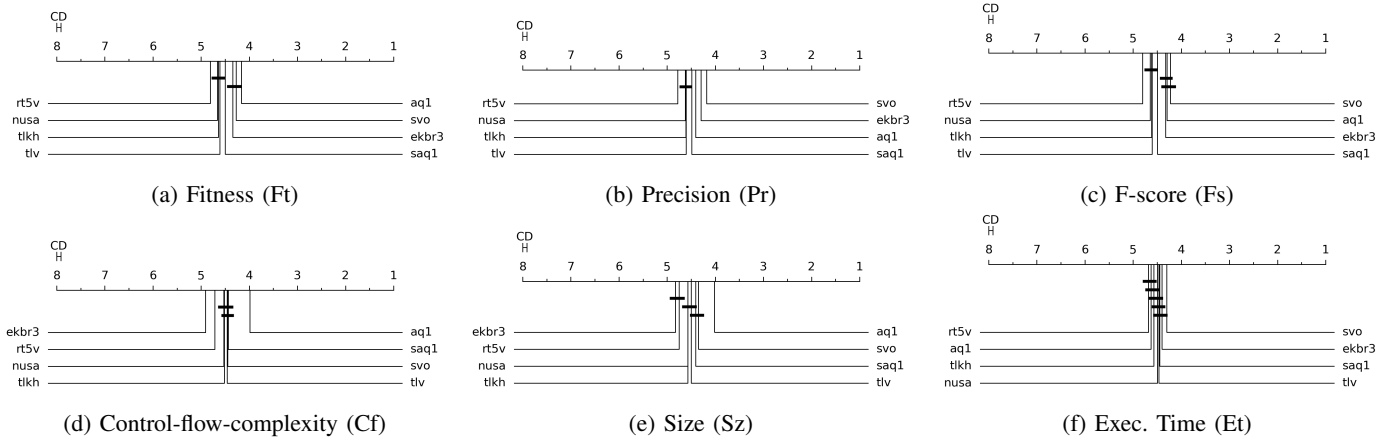


Fig. 3: CD diagrams showing meta-features' contributions to three selected metrics across all evaluated PD algorithms.

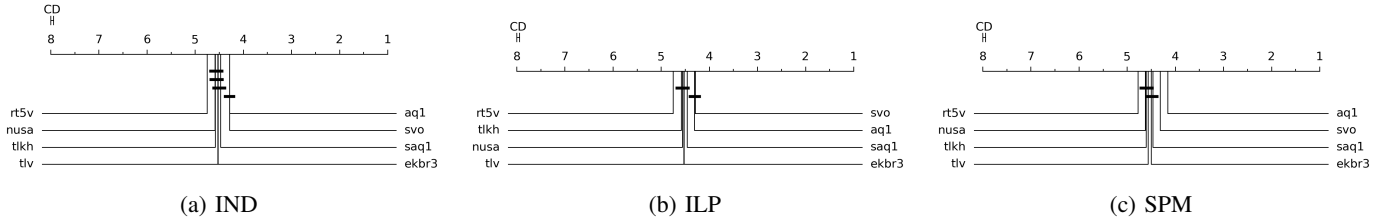


Fig. 4: CD diagrams showing significant meta-features' contributions per PD algorithm across all metrics

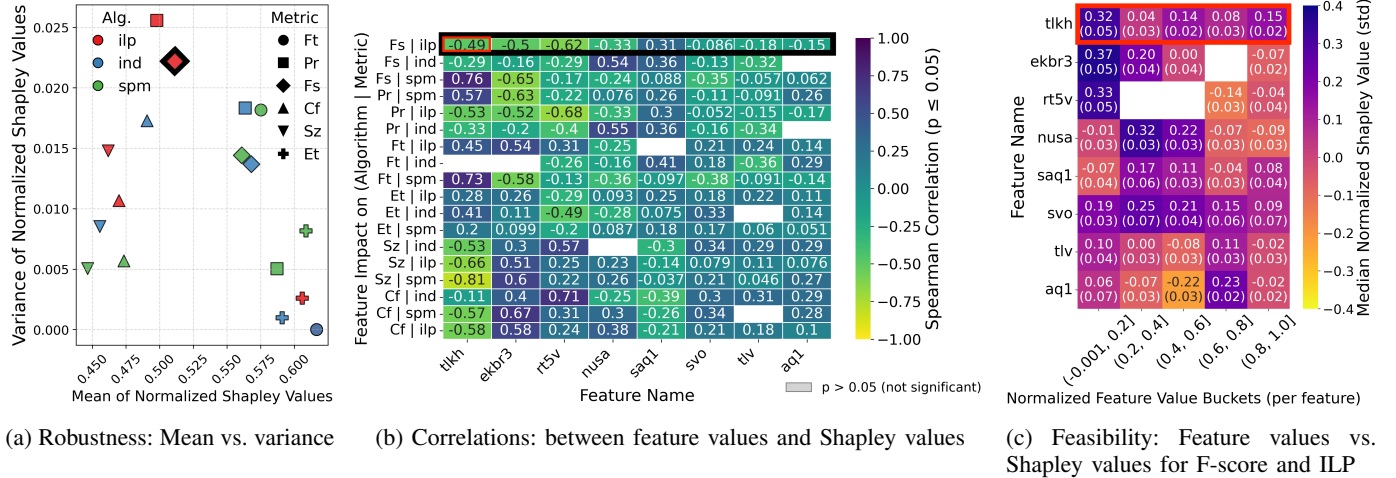


Fig. 5: *Robustness, correlation, and feasibility* of feature values and Shapley values for process discovery algorithms

models, as these features make it difficult for the miner to find a coherent solution. While previously discussed features have clear impact trends on metrics, the concentration of medium-length traces (*tikh*) shows a negative correlation with all discussed metrics. I.e., if most traces are of similar length with few exceptions, it can lead to overgeneralization, where *ILP* produces a simpler model that fails to capture the full range of process behavior, thereby reducing its *Pr* and *Fs*.

Feasibility. To explore the source of observed correlations and each algorithm's ability to handle different feature profiles, we use heatmaps of bucketed feature intervals versus normalized Shapley values. Figure 5c shows missing *Fs* values, indicating configurations where *ILP* failed to produce a model.

Feasibility issues arise for moderate-to-low *rt5v* values, where moderate variability and structural ambiguity increase search complexity. In contrast, extreme *rt5v* values often imply a clearer structure, making optimization more likely to succeed. Considering occasional unsoundness and resource constraints, Table I reports the percentage of event logs for which each algorithm produced viable models. Despite *IND*'s soundness guarantees[25], its feasibility was limited: Only 39% of logs completed within a five-minute timeout. This reflects *IND*'s vulnerability to time constraints when processing large or structurally complex logs. *SPM* achieves (76%) feasibility using a heuristic approach that bypasses strict constraints, but its sensitivity to high variability can cause failures by priori-

tizing precision over model structure. These results highlight a core challenge in process discovery: Algorithm feasibility and robustness vary with log characteristics and computational limits. *ILP* struggles with infeasibility, *IND* with timeouts, and *SPM* with structural trade-offs, underscoring the need to match algorithms to log characteristics and resource constraints.

C. RQ3: How can insights about the impact of EL characteristics support process miners on a specific downstream task?

In Table II we show the correlations between feature values and their impact on results according to their Shapley values. A positive arrow (\uparrow or \Uparrow) indicates that logs with larger feature values tend to co-occur with larger Shapley values. This means that when a feature value was set higher (lower), it systematically appears as more (or less) impactful in the average marginal contribution. However, this does not imply a causal relationship, as Shapley analysis is associational. We assess the utility of SHAining in PD from three perspectives: algorithm design, algorithm evaluation, and insights to build comprehensive algorithm libraries.

Algorithm Design. *Et* of *IND* is heavily influenced by the event log’s characteristics. High *Et* correlates with a greater concentration of medium-length traces (high *tlkh*) and an unbalanced distribution of variant occurrences (high *svo*). Conversely, logs with fewer rare variants (low *rt5v*) tend to reduce *Et* due to *IND*’s recursive, decomposing nature. Log characteristics that increase the depth of the recursion tree (*tlkh*), the number of branches (*rt5v*), or the difficulty of finding a suitable split (*svo*) increase computational complexity and thus *Et*. These findings highlight *IND*’s sensitivity to these features. This first invites algorithm developers to design an algorithm more robust to changes in those features regarding *Et*. Nevertheless, *ILP* and *SPM* are already more robust, exhibiting only an insignificant ($|r| < 0.1$) to low ($0.1 < |r| \leq 0.3$) correlation between these same log characteristics and *Et*, thus the suggested improvement has been achieved. On the other hand, we find a medium ($0.3 < |r| \leq 0.5$) to strong ($|r| > 0.5$) correlation between any other metric and at least one feature. While trade-offs between process discovery metrics [22] prevent one single universally robust algorithm, our analysis identifies opportunities. These insights could guide the development of new algorithms or optimizations designed to be more robust to changes in specific log characteristics, specifically against e.g., low *tlkh* and high *ekbr3* for *Sz*.

Algorithm Evaluation. When evaluating a process discovery algorithm, we can leverage a feature-dropout mechanism with ablations to analyse robustness, as in our framework. Stress tests using varying or constant event log characteristics offer transparent insight into the power and limitations of each algorithm. Observing the columns for *Fs*, between the algorithms, we can discover which feature variations they are the most robust against. Although for *SPM*, the concentration of medium-length traces (*tlkh*) and the entropy of triplets of activities (*ekbr3*) strongly impact the *Fs* results, these are less vulnerable to variations in the quantity (*nusa*) and the 25th

percentile of starting activities (*saq1*), because the algorithm’s frequency-based filtering and loop/concurrency detection directly rely on the variety of whole traces (*tlkh*, *ekbr3*), but less so on single starting activities (*nusa*, *saq1*). In contrast, the *Fs* results of *IND* are most strongly and negatively impacted by a smaller set of starting activities (low *nusa*). With only signs of insignificant to weak correlations for *tlkh* and *ekbr3*. While *ILP* depicts similar strongly correlated features as both *IND* and *SPM* for *Fs*, results underline the negative impact on *Fs* by high *ekbr3* and a high amount of rare traces (high *rt5v*), which introduce noise and complexity that make it difficult for the algorithm to create a precise and generalizable model. Likewise, a high *tlkh*, high *nusa*, or a low *saq1* can also lead to overly simplified or imprecise models, further lowering *Fs*. A novel evaluation framework could apply a masking on certain log features (e.g., collapsing of rare activity pairs, smoothing trace length variance) to test their effect on results quality.

Algorithm Libraries. Taking into account the metrics’ trade-offs [22] and quantifying the impact of event log characteristics, presented in our work, process miners can build broad libraries of PD algorithms, aiming at the best results for diverse incoming event logs. For example, consider how *IND* and *SPM* complement each other in *Fs* robustness, because their vulnerabilities are different, including disjunct strong correlation feature sets $((nusa, saq1, tlkv) \cap (tlkh, ekbr3, svo) = \emptyset)$. Similarly, *ILP* is a good addition to the (*IND*, *SPM*) library, considering that it has a reduced number of features (*tlkh*, *ekbr3*, *nusa*), which present vulnerabilities in comparison to *IND* and *SPM*. Similar arguments can be made in favor of *SPM* for *Ft* and *Pr*. Nevertheless, especially for *Sz*, an algorithm that shows robustness (insignificant to weak correlations across all features) or vulnerability towards variations in features other than (*tlkh*, *ekbr3*), would be a great addition to this library.

VI. CONCLUSION

PM algorithm evaluations often lack a systematic protocol to understand how event log characteristics influence model outcomes. We propose SHAining, a functional Shapley analysis that assesses how meta-level configuration parameters impact the algorithm metrics by varying the input event log characteristics, instead of fixed input event logs as in prior work. Our method follows four steps: (i) forming feature coalitions; (ii) sampling event logs conditioned on these coalitions; (iii) applying a PM algorithm on these logs; and (iv) performing Shapley value analysis to quantify feature contribution on evaluation metrics.

A large-scale evaluation across comprehensive event logs reveals how feature combinations affect different metrics in Process Discovery (PD). We identify the most impactful features across metrics and algorithms, uncover correlations between feature values and their contributions, and provide practical insights for PD algorithm design, evaluation, and libraries. Overall, SHAining enables understanding of how algorithm assumptions interact with varying event log characteristics, which is critical for PM where generalizability and model structure depend strongly on input characteristics.

TABLE I: Feasible Logs by Miner.

PD Algorithm	Feasible Logs [%]
IND	39
ILP	59
SPM	76
Overlap	39

TABLE II: Significant correlations (arrows denote direction; gray if $|r| < 0.1$, bold if $|r| > 0.3$, double arrow if $|r| > 0.5$)

Feature	ILP						IND						SPM					
tlkh	↑Ft	↓Pr	↓Fs	↓Cf	↓Sz	↑Et	-	↓Pr	↓Fs	↓Cf	↓Sz	↑Et	↑Ft	↑Pr	↑Fs	↓Cf	↓Sz	↑Et
ekbr3	↑Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	-	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	↓Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et
rt5v	↑Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	↓Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	↓Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et
nusa	↓Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	↓Ft	↑Pr	↑Fs	↓Cf	-	↓Et	↓Ft	↑Pr	↓Fs	↑Cf	↑Sz	↑Et
saq1	-	↑Pr	↑Fs	↓Cf	↓Sz	↑Et	↑Ft	↑Pr	↑Fs	↓Cf	↓Sz	↑Et	↓Ft	↑Pr	↑Fs	↓Cf	↓Sz	↑Et
svo	↑Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	↑Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	↓Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et
tlv	↑Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	↓Ft	↓Pr	↓Fs	↑Cf	↑Sz	-	↓Ft	↓Pr	↓Fs	-	↑Sz	↑Et
aq1	↑Ft	↓Pr	↓Fs	↑Cf	↑Sz	↑Et	↑Ft	-	-	↑Cf	↑Sz	↑Et	↓Ft	↑Pr	↑Fs	↑Cf	↑Sz	↑Et

Threats to Validity. While our FEEED feature selection provides broad coverage, other relevant features may exist. To ensure more reliable analysis, we applied an intercorrelation-based feature selection procedure (cf. Section V) because Shapley value analysis relies on low feature intercorrelation. Our evaluation may also be weighted towards chosen metric types and PD algorithms, which may limit the generalizability of our results by not covering all possible approaches. Nevertheless, our modular framework allows easy incorporation of additional features, algorithms, and metrics in future work.

Future Work. Future investigations should address these limitations by expanding the feature space beyond FEEED, exploring additional methods to further reduce feature intercorrelation effects on Shapley values, and incorporating a broader and more diverse set of process discovery algorithms.

We plan to explore computational optimizations, including refined sampling, to speed up our pipeline. Moreover, we'll address generator bias by testing alternative generators and validating our synthetic data, as we secure diverse real-world logs. Finally, examining industrial characteristics will further enhance our approach's practical relevance.

REFERENCES

- [1] J.-R. Rehse, S. Leemans, P. Fetteke, and J. M. E. van der Werf, "On process discovery experimentation," *ACM TOSEM*, 2024.
- [2] J. M. E. van der Werf, A. Polyvyanyy, B. R. van Wensveen, M. Brinkhuis, and H. A. Reijers, "All that glitters is not gold: Towards process discovery techniques with guarantees," in *CAiSE*, pp. 141–157, Springer, 2021.
- [3] A. Maldonado, C. M. M. Frey, G. M. Tavares, N. Rehwald, and T. Seidl, "GEDI: generating event data with intentional features for benchmarking process mining," in *BPM*, pp. 221–237, Springer, 2024.
- [4] A. Burattin, B. Re, L. Rossi, and F. Tiezzi, "A purpose-guided log generation framework," in *BPM*, pp. 181–198, Springer, 2022.
- [5] T. Jouck and B. Depaire, "Generating artificial data for empirical analysis of control-flow discovery algorithms," *BISE*, pp. 695–712, Dec 2019.
- [6] J. Munoz-Gama, N. Martin, C. Fernandez-Llatas, O. A. Johnson, M. Sepúlveda, E. Helm, V. Galvez-Yanjari, E. Rojas, A. Martinez-Millana, *et al.*, "Process mining for healthcare: Characteristics and challenges," *JBIS*, vol. 127, p. 103994, 2022.
- [7] T.-H. Huang, T. Junied, M. Pegoraro, and W. M. van der Aalst, "Proreco: A process discovery recommender system," in *CAiSE*, pp. 93–101, Springer, 2024.
- [8] C. Schreiber, "Exploring the impact of process diversity on business process performance," in *ICPM DC*, pp. 17–18, CEUR-WS.org, 2021.
- [9] S. K. vanden Broucke, C. Delvaux, J. Freitas, T. Rogova, J. Vanthienen, and B. Baesens, "Uncovering the relationship between event log characteristics and process discovery techniques," in *BPM Workshops*, pp. 41–53, Springer, 2013.
- [10] L. S. Shapley, *A Value for N-Person Games*. RAND Corporation, 1952.
- [11] S. Lipovetsky, "Handbook of the shapley value," *Technometrics*, vol. 62, no. 2, pp. 1–280, 2020.
- [12] J. Carmona, M. de Leoni, B. Depaire, and T. Jouck, "Summary of the process discovery contest 2016," in *ICPM*, Springer, 2017.
- [13] A. Stevens, J. De Smedt, and J. Peepkorn, "Quantifying explainability in outcome-oriented predictive process monitoring," in *ICPM*, pp. 194–206, Springer, 2021.
- [14] M. Pishgar, S. Harford, J. Theis, W. Galanter, J. M. Rodríguez-Fernández, L. Chaisson, Y. Zhang, A. Trotter, K. M. Kochendorfer, *et al.*, "A process mining-deep learning approach to predict survival in a cohort of hospitalized covid-19 patients," *BMC Medical Informatics and Decision Making*, p. 194, 2022.
- [15] T. Heskes, E. Sijben, I. G. Bucur, and T. Claassen, "Causal shapley values: Exploiting causal knowledge to explain individual predictions of complex models," in *NeurIPS*, 2020.
- [16] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *NeurIPS*, vol. 30, Curran Associates, Inc., 2017.
- [17] J. Marques-Silva and X. Huang, "Explainability is Not a game," *Commun. ACM*, vol. 67, no. 7, pp. 66–75, 2024.
- [18] J. Marques-Silva, "Logic-based explainability: Past, present and future," in *12th International Symposium, ISO LA 2024, October 27–31, 2024, Proceedings, Part IV*, p. 181–204, Springer-Verlag, 2024.
- [19] Y. Kwon, M. A. Rivas, and J. Zou, "Efficient computation and analysis of distributional shapley values," in *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, pp. 793–801, 2021.
- [20] A. Augusto, R. Conforti, M. Dumas, M. La Rosa, F. M. Maggi, A. Marrella, M. Mecella, and A. Soo, "Automated discovery of process models from event logs: Review and benchmark," *IEEE TKDE*, vol. 31, no. 4, pp. 686–705, 2018.
- [21] K. Andree, M. Hoang, F. Dannenberg, I. Weber, and L. Pufahl, "Discovery of workflow patterns - a comparison of process discovery algorithms," in *CoopIS*, pp. 257–274, Springer, 2024.
- [22] G. Janssenswillen, N. Donders, T. Jouck, and B. Depaire, "A comparative study of existing quality measures for process discovery," *Information Systems*, vol. 71, pp. 1–15, 2017.
- [23] A. Maldonado, G. M. Tavares, R. S. Oyamada, P. Ceravolo, and T. Seidl, "FEEED: feature extraction from event data," in *ICPM Demos*, 2023.
- [24] C. Back, S. Debois, and T. Slaats, "Entropy as a measure of log variability," *Journal on Data Semantics*, vol. 8, p. 129–156, June 2019.
- [25] S. J. Leemans, D. Fahland, and W. M. Van Der Aalst, "Discovering block-structured process models from event logs-a constructive approach," in *Petri Nets*, pp. 311–329, Springer, 2013.
- [26] S. J. van Zelst, B. F. van Dongen, W. M. van der Aalst, and H. Verbeek, "Discovering workflow nets using integer linear programming," *Computing*, vol. 100, pp. 529–556, 2018.
- [27] A. Augusto, R. Conforti, M. Dumas, M. La Rosa, and A. Polyvyanyy, "Split miner: automated discovery of accurate and simple business process models from event logs," *KAIS*, vol. 59, pp. 251–284, 2019.
- [28] A. Augusto, J. Mendling, M. Vidgof, and B. Wurm, "The connection between process complexity of event sequences and models discovered by process mining," *Information Sciences*, vol. 598, pp. 196–215, 2022.
- [29] B. Dongen, J. Carmona, and T. Chatain, "Alignment-based quality metrics in conformance checking," *EMISA Forum*, pp. 77–80, 2016.
- [30] A. Adriansyah, J. Munoz-Gama, J. Carmona, B. F. Van Dongen, and W. M. Van Der Aalst, "Measuring precision of modeled behavior," *ISEB*, vol. 13, pp. 37–67, 2015.
- [31] A. Adriansyah, B. F. van Dongen, and W. M. van der Aalst, "Conformance checking using cost-based fitness analysis," in *EDOC*, 2011.

- [32] J. Mendling, *Metrics for process models: empirical foundations of verification, error prediction, and guidelines for correctness*. Springer, 2008.
- [33] W. M. Van der Aalst, “Verification of workflow nets,” in *Petri Nets*, pp. 407–426, Springer, 1997.
- [34] S. Herbold, “Autorank: A python package for automated ranking of classifiers,” *Journal of Open Source Software*, p. 2173, 2020.