# A Conformance Checking-based Approach for Drift Detection in Business Processes

Víctor Gallego-Fontenla, Juan C. Vidal, and Manuel Lama
Centro Singular de Investigación en Tecnoloxías Intelixentes (CiTIUS),
Universidade de Santiago de Compostela, Galicia, Spain.
E-mail: {victorjose.gallego, juan.vidal, manuel.lama}@usc.es

*Abstract*— **Real life business processes change over time, in both planned and unexpected ways. The detection of these changes is crucial for organizations to ensure that the expected and the real behavior are as similar as possible. These changes over time are called concept drift and its detection is a big challenge in process mining since the inherent complexity of the data makes difficult distinguishing between a change and an anomalous execution. In this paper, we present *C2D2* (Conformance Checking-based Drift Detection), a new approach to detect sudden control-flow changes in the process models from event traces. *C2D2* combines discovery techniques with conformance checking methods to perform an offline detection. Our approach has been validated with a synthetic benchmarking dataset formed by 68 logs, showing an improvement in the accuracy ($F_{score}$) while maintaining a minimum delay in the drift detection.**

## I. Introduction

Real life processes are not immutable. Instead, they evolve to adapt to changes in their context, as new regulations or new consumption patterns. Changes can be planned by the organisation, but also happen unexpectedly. In the first case, the impact on the process can be computed and minimized. But in the second case, it may lead to wrong decisions because of outdated information. Thus, organizations should put in place prevention measures to detect when something is running differently from planned to reduce this negative impact. These unforeseen changes over time are known as concept drifts, which is one of the challenges presented in the *Process Mining Manifesto* [1].

Changes can be classified based on their distribution over time [2]: (i) sudden drifts (Figure 1a), which means that the new concept replaces totally the previous one; (ii) gradual drifts (Figure 1b), where the new and the old concepts coexist for some time; and (iii) incremental drifts (Figure 1c), when the transition from the oldest concept to the newest one passes through some intermediate states that are, usually, some kind of combination from both. Furthermore, when changes can be repeated over time, periodically switching between concepts, the change is classified as a recurrent drift (Figure 1d). In this paper, we focus on sudden drift detection.

In addition, based on how data are processed [3], concept drift can be: (i) offline, when change detection is made *post–mortem*, being all data available from the beginning, and (ii) online, when change detection is made *on-the-fly*, and new data

are processed just when it is generated. In this paper we focus on offline concept drift, which additionally faces two challenges: a) the inherent complexity of process models, that can contain and combine different structures such as sequences, loops, parallel branches and choices; and b) the distinction between a change and an outlier, which is not always clear and may depend on the application domain.

Although, some authors have proposed different approaches for concept drift detection in process mining [3]–[18], identifying all the possible change patterns [19] with a short delay, allowing organizations to know exactly when the change took place and helping in the identification of the reasons that caused the change, is still a challenge. In addition, many approaches are unable to detect all change patterns, being this essential to reduce the cases when a change stays ignored to the organization. Another issue in some proposals is their high dependence on the end user, who is required to have some *a-priori* knowledge of the process structure or skills to identify accurately the drift within a set of possibilities.

In order to reduce the aforementioned issues, in this paper we present *C2D2*, a novel and fully automatic approach based on discovery and conformance checking techniques for offline detection of sudden concept drifts in the control-flow of process models. The method starts by defining a reference window, that will serve as a ground truth. Then traces are processed by a discovery algorithm to extract the corresponding process model. To process the remaining traces, the window is slided over the log, updating conformance metrics related to that process model. With these conformance values a regression is computed, and when the measurements decrease significantly a drift is detected. The underlying idea is that the value from the conformance metrics computed over the reference model and the new traces should decrease when the latter comes from a modified process, being this enough to determine if a change exists or not with a low delay. Specifically, the main contribution of this paper is the use of conformance metrics, in particular, fitness and precision, to detect changes in processes, which is a novel and unexplored approach so far. Namely, we propose the use of fitness metrics to detect changes that include traces with behavior not supported by the current process model, and the use of precision metrics to detect changes that implies behavior from the model disappearing from the real executions.
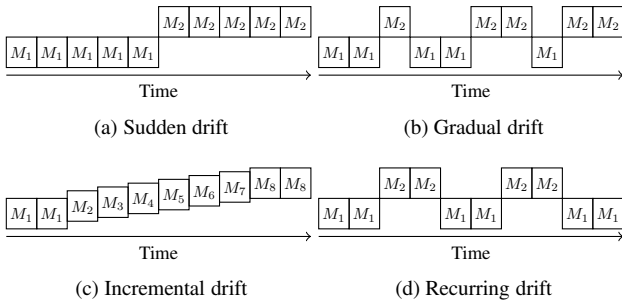
Figure 1. Types of concept drift based on their occurrence over time.

*C2D2* has been tested using a dataset with 68 synthetic event logs. The results have been compared with the ones obtained by the methods available in the state of the art. *C2D2* has proved to be better at the accuracy level, getting better $F_{score}$ for all the event logs. In addition, *C2D2* gets very low delays, identifying changes closer to the point in which they really happened. Getting good values for both metrics is important for organizations for minimizing the number of unidentified changes and for reacting as soon as possible to those changes.

The rest of this paper is structured as follows. In Section II we analyze the main approaches to concept drift analysis. In Section III we define a set of terms necessary to understand correctly our approach. In Section IV we detail our method for offline control-flow process concept drift detection. In Section V we present the experimentation performed to validate our approach and how it outperforms the main algorithms from the literature. Finally, in Section VI we present our conclusions and outline our future work.

## II. RELATED WORK

Although process mining is a rather active research field, concept drift analysis has not received much attention until recently. It is worth noting that, although the method proposed in this paper focuses on offline detection, online approaches are also considered in the following analysis, because they can be easily adapted to detect this type of changes by simulating an online environment from the complete event log.

In [4], authors propose a method for online concept drift detection using a polyhedron-based log representation. Then, they monitor the probability that a trace falls into that polyhedron using the ADWIN algorithm [20]. The main drawback of this approach is that it can only detect the presence of a change, but it does not give any information about when it happened.

Online detection is also addressed in [5], where authors discover a probabilistic process model that, given an activity, assigns probability to every possible successors, and check how this probabilities evolve throughout the complete log using statistical hypothesis tests. Although the method identifies drifts in most cases, small changes in less likely activities generate changes in the probabilities that can stay undetected.

In [6] authors propose an online approach based on the extraction of histograms from traces, and then use a clustering algorithm to generate groups of similar traces. A change is triggered when a new cluster appears. An important drawback

of this approach is that events order is not accounted for. Thus, it can only detect the addition or removal of new activities, but not the changes in the precedence relations between them.

In [3], authors use a fixed-size window over some features extracted from the follows/precedes relations present in traces, and statistical hypothesis tests to evaluate whether these features have changed significantly. The weak point of this method is that it requires a lot of interaction from the user, including previous knowledge of the process model and the areas where the changes can be located. An extension of this work has been proposed in [8], where authors implement a recursive bisectioning approach. Specifically, they take the traces that are involved in a drift detection and recursively split them in halves, with the aim of automatically localizing the change. A drawback of this approach is that it still requires the user to know the possible changes in order to obtain good results. A similar solution is presented in [9], where authors propose the usage of event class correlation as feature, and apply statistical hypothesis tests to detect changes. However, it fails in detecting some change patterns such as the changes in the execution order of activities.

Another approach followed by some authors is the usage of clustering techniques in order to detect the drift. In [11], authors cluster traces using the distance between pairs of activities. However, this approach does not support models with loops. Moreover, the distance can ignore certain change patterns depending on how many activities are affected by the change. In [12], authors extend a trace clustering algorithm [21] adding a time dimension in order to force clusters to include only consecutive traces, and thus be able to detect changes. Their approach highly depends on the number of clusters, fixed by the user, and only obtain good results when the number of clusters is equal to the number of changes. In [13], authors use a Markov clustering algorithm over different time windows in order to detect changes, but the approach does not focus on the control-flow perspective. Instead, multiple viewpoints of the process are taken into account simultaneously, mixing control-flow changes with behavioural and resources changes.

Another interesting approach, called *ProcessDrift*, is proposed in [14], where authors transform traces into *partial-ordered-runs* and then apply a statistical hypothesis test over two windows (one for reference and one for detection) in order to detect changes. The main drawback of this approach lies in its sensitivity to changes in the frequencies of certain relations present in the log, that may lead to false positives in the detection. A related method is presented in [18], where authors focus on detecting the change at the event level instead of at trace level. Specifically, they extract the $\alpha^+$ relations from two consecutive adaptive windows of events, and then, applying a statistical test, namely the G-test, compare the relations distribution of these two windows. This allows the detection even with unfinished executions, and reduce the detection delay. The drawback of this approach is that it requires high amounts of traces to be able to detect changes, being possible to ignore them when they are close from each other.

In [15], [16] authors apply graph metrics to detect changes.

In [15], authors compare the eigenvectors and the eigenvalues of undirected weighted graphs representing the log at different instants. In this graph, each vertex represents a trace. The edges weight is the similarity between the vertex (traces) it connects. However, this method needs a huge amount of traces, being unable to detect changes in logs with less than 2,000 traces. In [16] authors compare models over time using graph features, such as the node degree, the graph density or the occurrence of nodes and edges. However, this approach does not perform well in processes with loops.

Finally, in [17] authors present *TPCDD*, a method that transforms the event log into a relation matrix using direct succession and weak order relations, where each column represents a trace and each row a relation. Then, based on the trend of these relations generate candidate drift points. These points are clustered using DBSCAN, in order to group candidates that belong to the same drift point. This approach relies heavily in the user defining a correct radius for the DBSCAN algorithm, potentially getting a high number of false positives when it is too low and a high number of false negatives when it is too high.

With *C2D2* we take the aforementioned issues and try to minimize them in order to improve the results of the process drift detection. The method removes any user interaction in the drift detection, requiring only a reference window size to be specified. Moreover, the method is able to detect all change patterns independently of the process structure. Furthermore, the method is designed to identify drifts with low delay, minimizing the detection of false negatives and positives.

## III. PRELIMINARIES

Below we present some concepts needed to understand the proposed method. The method takes an event log resulting of the executions of a process and tries to detect the changes in the execution of that process over time.

**Definition 1** (**Event**). An event $\varepsilon$ represents the execution of the activity $\alpha$ in the context of a process. Events have some mandatory attributes such as the activity, the execution case or the execution timestamp. They can also have optional attributes, such as the resource that performed the activity, the variables that were modified or the location.

**Definition 2** (**Trace**). A trace is an ordered sequence of events $\tau = \langle \varepsilon_1, ..., \varepsilon_n \rangle$ where every event belongs to the same execution case.

**Definition 3** (**Log**). A log is defined as an ordered collection of traces $L = \langle \tau_1, ..., \tau_n \rangle$ where each trace represents one execution of the process. The size of the log, denoted as $|L|$, represents the number of traces in that log.

A process model is a grah that describes the log behavior, that is, that allows to replay the log traces. A process model contains a representation of the coordination between the process activities, through sequences, parallels, choices and so on. In this paper we formalize process models using Petri nets.

**Definition 4** (**Petri net**). A Petri net is a tuple $N = (P, T, F)$, where:
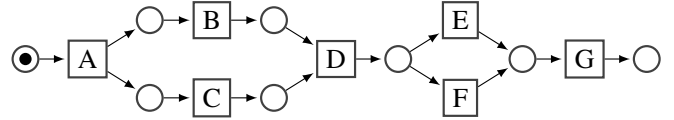


Figure 2. Example of a Petri net representing a process model.

- $P$ is a finite set of places
- $T$ is a finite set of transitions;
- $P \cap T = \emptyset$; and
- $F \subseteq (P \times T) \cup (T \times P)$ is a set of directed arcs.

Given $x \in T \cup P$, the set $^\bullet x = \{y \mid (y, x) \in F\}$ is the set of inputs of $x$, and $x^\bullet = \{y \mid (x, y) \in F\}$ the set of outputs of $x$. Given a Petri net $N = (P, T, F)$, a marking of $N$ is a mapping $M : P \to \mathbb{N}$, where $\mathbb{N}$ is the number of tokens in the place.

Processes usually have a unique start place $s \in P$ which has no inputs ($^\bullet s = \emptyset$) and a unique end place $f \in P$ which has no outputs ($f^\bullet = \emptyset$). The initial marking of the Petri net $M_0$ contains only the initial place $M_0(s) = 1 \wedge \forall q \neq s \in P : M_0(q) = 0$. For a transition $t$ to be fired, all its input places must contain at least one token ($\forall p \in {}^\bullet t : M(p) \geq 1$). When $t$ is executed, it consumes a token from each of its inputs and puts a token in every of its outputs. Petri nets can be depicted as bipartite graphs, being transitions represented as rectangles and places as circles. A black bullet into a place represents a token.

Figure 2 shows a Petri net example. In this example, the process is conformed by the activities $A$, $B$, $C$, $D$, $E$, $F$ and $G$. In real executions, $A$ must be executed first. Then, $B$ and $C$ can be executed in any order. After this two activities are finished, $D$ is executed. Then, exclusively one of $E$ and $F$ must be executed. Finally, $G$ is executed and the process execution finishes.

The quality of a process model $N$ with respect to a log $L$ can be estimated through some well established metrics such as fitness and precision.

**Definition 5** (**Fitness metric**). Given a log $L$ and a process model $N$, the fitness can be defined as a function $\gamma : L \times N \to \mathbb{R}$ which measures how well the log traces can be executed by the model.

**Definition 6** (**Precision metric**). Given a log $L$ and a process model $N$, a precision measure can be defined as a function $\rho : L \times N \to \mathbb{R}$ which measures how much additional behaviour is allowed by the model that is never observed in the log.

One of the objectives in this paper is the identification of changes in the process structure over time. This is done exploring the traces generated by the process and assessing if the more recent traces are product of the same process model. The concept of a sliding window captures this latter set of traces.

**Definition 7** (**Sliding window**). Given a log $L$ and an integer $n \leq |L|$, a sliding window of size $n$ over the log $L$ can be defined as the sublog that at instant $i$ contains the last $n$ traces, denoted by $\omega_i = [\tau_i, ..., \tau_{i+n}]$. When a new trace is read from the log, $i$ is incremented, so the oldest trace from the window
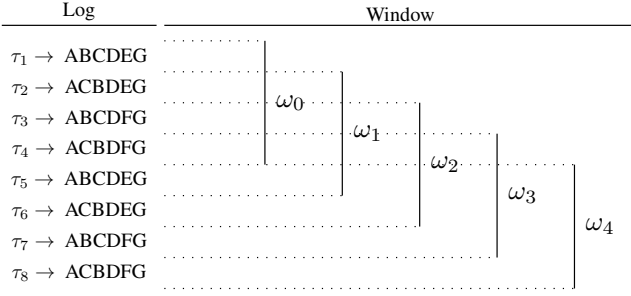
Figure 3. Example of the sliding window behaviour.

is forgotten, and the new one is added at the end of the window ($\omega_{i+1} = [\tau_{i+1}, ..., \tau_{i+n+1}]$).

Figure 3 shows an example log and the behaviour of a four-sized sliding window over it. At instant $i = 0$, the sliding window $\omega_0$ contains traces $\tau_1$ to $\tau_4$. When a new trace is read from the log, $i$ is incremented, so the oldest trace in the window ($\tau_1$) is forgotten and the new trace ($\tau_5$) is added to the window. This behaviour continues until the full log has been read.

The structural evolution of the process over time, e.g., to adapt to new context circunstances or organization needs, is known as control-flow drift.

**Definition 8 (Control-flow drift).** Let $N_i = (P, T, F)$ and $N_j = (P^*, T^*, F^*)$ be two process models at instants $i$ and $j = i + 1$, respectively. We say a control-flow drift exists when any of the following conditions is satisfied:

- $T \neq T^*$, which means that there are different activities in both models,
- $F \neq F^*$, which means that the relations between activities have changed.

## IV. $C2D2$: Conformance Checking-based Concept Drift Detection

The pseudocode of the proposed algorithm is listed in Algorithm 1. Drift detection is based on a fixed-size sliding-window (definition 7), and thus the window size is the only parameter required to be specified by the user, aside from the event log. The rest of this paper assumes a window defined as a number of traces, but the method can be easily modified to allow windows based on a time interval.

The algorithm starts by initializing a list $D$ of trace indices that produce a drift detection (line 2), and the initial window index $i = 0$ (line 3). Then, for this initial window, two lists of fitness and precision measurements $\Gamma$ and P (lines 5 and 6), and two lists of drift candidates $D^\Gamma$ and $D^P$ (lines 7 and 8), one for each metric, are initialized. This two lists $D^\Gamma$ and $D^P$ contain a boolean for each processed window, indicating whether the window has been marked as a drift candidate or not.

A reference model $N$ is discovered from the window content $\omega_i$ (line 10). $C2D2$ is not tied to any discovery algorithm, although some features can affect its performance. First, the computational complexity since the discovery algorithm can be executed multiple times, one for each detected drift. Second,

discovery algorithms that maximize fitness are more sensitive to small variations in traces, enhancing the subsequent detection of changes by traces replayability. This initialization phase is executed on the first trace window and when a new drift is detected by the algorithm.

The detection step comprises lines 11-17. As a first step, the fitness and precision metrics, $\gamma$ and $\rho$, are computed and appended to the respective list of measurements $\Gamma$ and P (lines 12 and 13). Then, the window is classified as a drift candidate for each metric using the method described in Section IV-A and the result of this classification (either true or false) is stored in the corresponding list of drift candidates $D^\Gamma$ and $D^P$ (lines 14 and 15). Finally, the window index is incremented by one (i.e., the window slides one position), reading a new trace from the log (line 16). A change is confirmed only if it persists in time. That means, the last $n$ windows has been classified as drift candidates either for fitness (lines 32-36) or precision (lines 37-41). This allows the method to prevent false positives due to the existence of temporal falls in the metrics caused by outlier traces.

Once drift is confirmed, index $i$ is added to the list of confirmed drifts $D$ (line 19), and the algorithm loops back to the initialization phase. Otherwise, detection step starts again with the new window, repeating until a change is detected or the log is completely analyzed.

### A. Drift Detection

Fitness and precision separately can not detect all possible changes, but a combination of both can. This is illustrated in the example depicted in Figure 4. Figure 4a and Figure 4b present two models, henceforth $N_1$ and $N_2$. The difference between both models is that activities $B$ and $C$ that are in parallel in $N_1$, while they are in sequence in $N_2$. Let suppose that a process changes from the model $N_1$ to $N_2$ at instant $i = 8$, which log is represented in Figure 4c, henceforth denoted as $L_1$. Traces $\tau_1$ to $\tau_8$ correspond to the execution of $N_1$ and traces $\tau_9$ to $\tau_{16}$ correspond to $N_2$. In $L_1$, the concurrent execution of activities $B$ and $C$ becomes a sequence from $\tau_9$ onwards. With this change there is not any trace that is not replayable by the model, so the fitness remains unaltered. However, no trace in the window contains the path $A \rightarrow C \rightarrow D$ from $\tau_9$ onwards, so the precision falls. This can be seen in Figure 4e, where precision falls because the model allows more behaviour than is present in the traces, but fitness remains unaltered.

Let also suppose a different change, from model $N_2$ to $N_1$, at the same time instant, which log is represented in Figure 4d, henceforth denoted as $L_2$. Traces $\tau_1$ to $\tau_8$ are generated by $N_2$ and traces $\tau_9$ to $\tau_{16}$ by $N_1$. In $L_2$, $B$ and $C$, originally in sequence, are in parallel from $\tau_9$ onwards. This change can not be detected using precision (the model does not generate more behavior than the present in the log), but it can be detected though fitness, since $\tau_{10}, \tau_{12}, \tau_{14}$ and $\tau_{16}$ can not be replayed by $N_2$. This situation is represented in Figure 4f, where precision remains unchanged, but fitness falls in the 7[th] iteration of the algorithm. This detection mechanism is listed in lines 24 to 30 of Algorithm 1 (function CONFIRMDRIFT). As a first step, a

**Algorithm 1** Conformance Checking-based Drift Detection

**Inputs:** an event log $L$ and a window size $n < |L|$
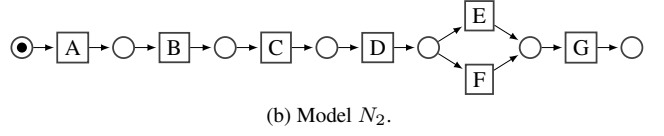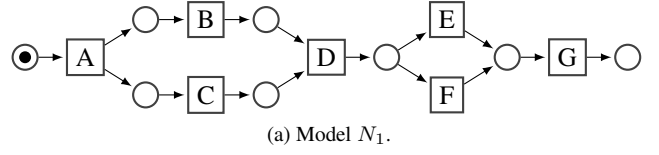**Outputs:** a list of trace indices causing drift D

```
1:  procedure CONCEPTDRIFTDETECTION(L, n)
2:      D ← [ ]
3:      i ← 0
4:      while i ≤ |L| − n do
5:          Γ  ← [ ] //fitness measures (Def. 5)
6:          P  ← [ ] //precision measures (Def. 6)
7:          D^Γ ← [ ] //drift candidates (fitness)
8:          D^P ← [ ] //drift candidates (precision)
9:          ω_i ← {L_i, ..., L_{i+n}}
10:         N  ← discover(ω_i)
11:         while (i ≤ (|L| − n)) ∧
                        ¬ CONFIRMDRIFT(n, D^Γ, D^P) do
12:             Γ  ← Γ :: γ(ω_i, M) //append current fitness
13:             P  ← P :: ρ(ω_i, M) //append current precision
14:             D^Γ ← D^Γ :: IDENTIFYDRIFT(n, Γ, D^Γ)
15:             D^P ← D^P :: IDENTIFYDRIFT(n, P, D^P)
16:             i  ← i + 1
17:         end while
18:         if CONFIRMDRIFT(n, D^Γ, D^P) then
19:             D ← D :: i
20:         end if
21:     end while
22:     return D
23: end procedure

24: function IDENTIFYDRIFT(n, data, D*)
25:     Υ  ← Θ([data_{|data|−(n/2)}, ..., data_{|data|}])
26:     m^< ← m(Υ) < 0 ∧ p(Υ) < 0.05
27:     m^> ← m(Υ) > 0 ∧ p(Υ) < 0.05
28:     m^= ← (¬ m^<) ∧ (¬ m^>)
29:     return (|data| > n/2) ∧ (m^< ∨ m^> ∨ (m^= ∧ (D*_{|D*|} = true)))
30: end function

31: function CONFIRMDRIFT(n, D^Γ, D^P)
32:     if {∀i ∈ [|D^Γ| − n, |D^Γ|] : D^Γ_i = true} then
33:         d^Γ ← true
34:     else
35:         d^Γ ← false
36:     end if
37:     if {∀i ∈ [|D^P| − n, |D^P|] : D^P_i = true} then
38:         d^P ← true
39:     else
40:         d^P ← false
41:     end if
42:     return (|D^Γ| ≥ n ∧ d^Γ) ∨ (|D^P| ≥ n ∧ d^P)
43: end function
```



(a) Model $N_1$.

(b) Model $N_2$.

| ID | Trace | ID | Trace |
|---|---|---|---|
| $\tau_1$ | A B C D E G | $\tau_1$ | A B C D E G |
| $\tau_2$ | A C B D E G | $\tau_2$ | A B C D F G |
| $\tau_3$ | A B C D F G | $\tau_3$ | A B C D E G |
| $\tau_4$ | A C B D F G | $\tau_4$ | A B C D F G |
| $\tau_5$ | A B C D E G | $\tau_5$ | A B C D E G |
| $\tau_6$ | A C B D E G | $\tau_6$ | A B C D F G |
| $\tau_7$ | A B C D F G | $\tau_7$ | A B C D E G |
| $\tau_8$ | A C B D F G | $\tau_8$ | A B C D F G |
| $\tau_9$ | A B C D E G | $\tau_9$ | A B C D E G |
| $\tau_{10}$ | A B C D F G | $\tau_{10}$ | A C B D E G |
| $\tau_{11}$ | A B C D E G | $\tau_{11}$ | A B C D F G |
| $\tau_{12}$ | A B C D F G | $\tau_{12}$ | A C B D F G |
| $\tau_{13}$ | A B C D E G | $\tau_{13}$ | A B C D E G |
| $\tau_{14}$ | A B C D F G | $\tau_{14}$ | A C B D E G |
| $\tau_{15}$ | A B C D E G | $\tau_{15}$ | A B C D F G |
| $\tau_{16}$ | A B C D F G | $\tau_{16}$ | A C B D F G |

(c) Log generated using $N_1$ as the initial process and $N_2$ as the modified one.

(d) Log generated using $N_2$ as the initial process and $N_1$ as the modified one.



(e) Fitness and precision evolution when using a window of size 4, log from Figure 4c and reference model $N_1$

(f) Fitness and precision evolution when using a window of size 4, log from Figure 4d and reference model $N_2$.

Figure 4. Measurements evolution for two logs with different changes.

linear regression is computed over the last $n/2$ measurements for both fitness and precision (line 25). A statistical test is performed to ensure that results are statistically significant, determining the probability of the regression slope being 0 with a $p_{value}$ of 0.05. When $H_0$ is rejected (i.e. $p(\Upsilon) < 0.05$) we asume that enough evidences exist to accept the slope value $m(\Upsilon)$. Otherwise, we can not assume that the slope value is different from 0. There are three possible situations:
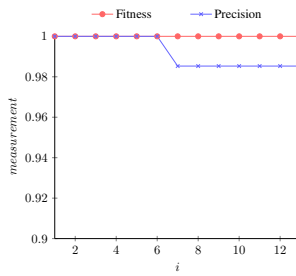
1) The regression slope is negative (line 26): metrics get lower values, so more traces are not replayable for fitness or, conversely, more paths of the model are not contained in traces for precision. Thus the window is marked as a drift candidate.
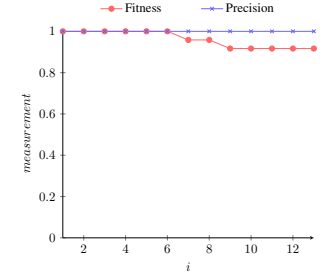
2) The regression slope is positive (line 27): metrics get higher values, since more traces are replayable for fitness or, conversely, more paths of the model are contained in traces for precision. Thus the window is marked as a drift candidate.

3) The regression slope is zero (line 28): no change in conformance metrics, i.e., the window does not present any drift. In this case, a drift can also be detected, but only if the previous window was marked as a drift candidate.

We can see the negative regression slopes in precision and fitness, for the aforementioned logs, in Figure 4e and Figure 4f,

(a) Log 1 (`pl`). Precision falls every time a change happens but fitness remain unaltered.



(b) Log 2 (`cb`). Fitness falls but precision remains unaltered in every change.

Figure 5. Evolution of fitness and precision metrics when computed with the last 100 traces on logs `pl` and `cb`.



| Log |
|---|
| A B C D E G |
| A C B D E G |
| A B C D E G |
| A C B D E G |
| A C B D E G |
| A B C D E G |
| A B C D E G |
| A C B D E G |
| A B C D E G |

$$OLP = \langle (A \to B)(A \to C)(B \to D)(C \to D)$$
$$(D \to E)(D \to F)(E \to G)(F \to G) \rangle$$

$$DFR = \langle (A \to B)(A \to C)(B \to C)(B \to D)$$
$$(C \to B)(C \to D)(D \to E)(E \to G) \rangle$$

$$PC = 1 - \frac{|\langle (D \to F)(F \to G) \rangle|}{8} = 0.75$$

Figure 6. $PC$ computation example.

respectively. A more complete example with several changes is depicted in Figure 5. Specifically, this example shows the drift detection using the logs `pl` and `cb`, that will be described in Section V-A, which contain a change every 250 traces. In the case of `pl` (Figure 5a), two fragments that are originaly executed in a concurrent form are transformed into a sequential execution, which should imply a reduction in precision but not in fitness. In the case of `cb` (Figure 5b), a fragment is transformed from mandatory to skipable, which should imply a reduction in fitness but not in precision. When the regression is computed with few data the slope can oscillate, producing false positives. To prevent this effect, the algorithm only checks the previous conditions if there are at least $n/2$ measurements available.
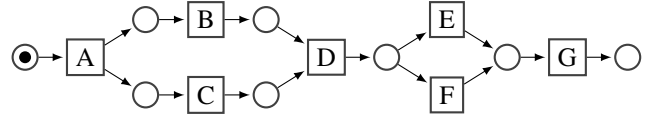
### B. Custom Fitness and Precision

Traditional fitness and precision metrics are designed to assess the global quality of a model. However, *C2D2* use them to detect structural changes in the execution of a process. The main intuitions behind this hypothesis are:

1) Fitness can measure how well the log can be replayed by the model. When new activities are added, removed or the precedence is modified, this metric will fall, indicating a change.
2) Precision can measure how well the model represents the behavior that is present in the log. When a path no longer exists, the precision values will fall, indicating a change.

Moreover, we also propose two simpler fitness and precision metrics aside from the well-established metrics from the state of the art [22], [23]. These two approaches have much lower computational complexity and were designed to detect changes in simple and noise-free logs. In the case of fitness, we use the percentage of replayable traces. This approach is not particularly useful for measuring the quality of a model, since it equally penalizes traces that do not fit the model and those that deviate slightly from it. Despite this, it can be used to estimate changes in fitness, since a change in the percentage of traces that can be replayed in the model always leads to a change in the metric value. For precision, the following approach is used:

$$PC = 1 - \frac{|OLP \setminus DFR|}{|OLP|} \tag{1}$$

where:

- A set of one-length paths ($OLP$) is extracted from the model. An $OLP$ is a pair of activities that are directly connected in the process model, without any other activity in between
- A set of directly-follows relations ($DFR$) is extracted from the log. A $DFR$ is a pair os activities that appear one after the other in the log, without any activity in between.

Eq. (1) does not measure the precision *per se*, but the change in the precision. The moment a $OLP$ stops appearing in the log is indicative that some path of the model has disappeared. The proposed approach returns $1$ when all the supported behavior of the model appear at least once in the log, and $0$ otherwise, i.e., when none of the behavior supported by the model appears in the log. The computation of this metric is illustrated with an example in Figure 6.

## V. Experimentation

Concept drift algorithms are assessed based on two quality measures: $F_{score}$ (2a), which is an accuracy metric computed as the harmonic mean between precision (2b) and recall (2c); and delay (henceforth $\Delta$), which is the distance between the point when the change really happened and when it is detected.

$$F_{score} = \frac{2 \times precision \times recall}{precision + recall} \tag{2a}$$

$$precision = \frac{TP}{TP + FP} \tag{2b}$$

$$recall = \frac{TP}{TP + FN} \tag{2c}$$

To classify the detected changes as *true positives* ($TP$), *false positives* ($FP$) or *false negatives* ($FN$), we use a threshold $\varepsilon$, that represents the error tolerance of the quality measures, and a neighborhood $\delta_i^\varepsilon$, defined as the interval between $i-\varepsilon$ and $i+\varepsilon$. Let a change happen at instant $i$. This change is classified as a $TP$ only when it is detected in $\delta_i^\varepsilon$. When no change is detected in $\delta_i^\varepsilon$ it is classified as $FN$. Finally, all changes detected in $\delta_i^\varepsilon$ where a previous change has been already detected is classified as a $FP$, as well as the ones detected outside any $\delta^\varepsilon$. Figure 7
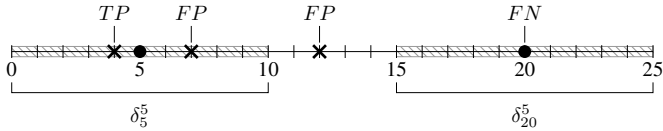
Figure 7. Change results classification in $TP$, $FP$ and $FN$. A dot represents a real change. A cross represents a detection. Shadowed in the neighborhood.

Table I
SIMPLE CHANGE PATTERNS FROM [19] APPLIED TO THE ORIGINAL MODEL.

| Code | Change pattern | Class |
|------|----------------|-------|
| cm | Move fragment into/out of conditional branch | I |
| cp | Duplicate fragment | I |
| pm | Move fragment into/out of parallel branch | I |
| re | Add/remove fragment | I |
| rp | Substitute fragment | I |
| sw | Swap two fragments | I |
| cb | Make fragment skippable/non-skippable | O |
| lp | Make fragments loopable/non-loopable | O |
| cd | Synchronize two fragments | R |
| cf | Make two fragments conditional/sequential | R |
| pl | Make two fragments parallel/sequential | R |

shows an example with two real changes ($d_5$, at instant 5, and $d_{20}$, at instant 20), and three detections, at instants 4 ($c_4$), 7 ($c_7$) and 12 ($c_{12}$), using a $\varepsilon = 5$. In this example, $c_4$ is classified as a $TP$, because it lies in the neighborhood of $d_5$; $c_7$ is classified as a $FP$, because, despite being in the neighborhood of $d_5$, another change has been detected previously; $c_{12}$ is classified too as a $FP$, in this case for being detected outside any neighborhood $\delta^5$; and finally, $d_{20}$ is classified as a $FN$ since no change is detected in its neighborhood.

Both the algorithm and the data used for the tests are available online[1].

*A. Validation Data*

To test the proposed approach, synthetic logs have been generated using the metodology, models and change patterns defined in [14]. This is the most accepted methodology for generating datasets when validating sudden concept drift detection algorithms in process mining. The generated dataset contains 68 logs: 17 with 2,500 traces, 17 with 5,000, 17 with 7,500 and 17 with 10,000. The original dataset from [14] contains 4 more logs (one for each of the sizes), but they have been discarded because its drifts (changing the frequency of the branches in a choice construct) are not control-flow drifts, but behavioural ones.

The original model used to generate the logs corresponds to a loan application process. The corresponding Petri net is depicted in Figure 9. To generate the 17 modified models, 11 simple change patterns from [19] are applied to the original process. The applied patterns are collected in Table I. These changes can imply an insertion (labeled as $I$), an optionalization (labeled as $O$) or a resequentialization (labeled as $R$). For

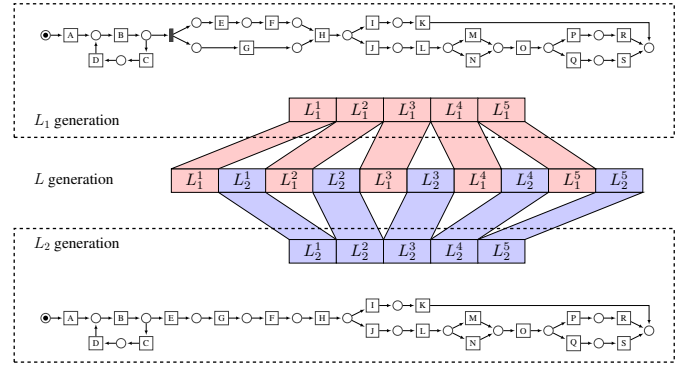[1]tec.citius.usc.es/concept-drift-api/swagger-ui.html



Figure 8. Log generation example.

each simple change pattern a different model is generated. The remaining 6 models are generated by applying a combination of simple change patterns, picking one change from each of the previously named classes.

Once all the models are available, the logs are generated simulating executions of those processes. The original log is then combined with the modified ones to generate logs with drifts. The final log is composed joining alternatively sublogs from both the original model and the modified ones. Each drifting log presents a change every $10\%$ of its final size. A log generation example is represented in Figure 8. Two logs ($L_1$ and $L_2$) with different models are split in 5 sublogs with equal sizes ($L_1^1$ to $L_1^5$ and $L_2^1$ to $L_2^5$). This sublogs are combined alternatively into a log $L$, with size $|L_1| + |L_2|$.

*B. Impact of Discovery Algorithm and Fitness and Precision Metrics in Detection*

In order to check the impact of the discovery algorithm and the conformance metrics in terms of both $F_{score}$, $\Delta$ and computation time, different configurations have been tested:

1) Discovery algorithms: Inductive Miner ($IM$) [24] and Heuristics Miner ($HM$) [25], which are two of the most used methods for discovering models from event logs. No algorithm based on evolutionary computation has been selected because it would increase the computational complexity significantly.
2) Fitness metrics: Alignment Based Fitness ($AF$) [22], Negative Event Recall ($NR$) [23] and the percentage of completely replayable traces ($RT$).
3) Precision metrics,: Advanced Behavioural Appropriateness ($ABA$) [22], Negative Event Precision ($NP$) [23] and precision change assessment ($PC$).

**Impact on $F_{score}$ and $\Delta$.** Table II shows the mean $F_{score}$ and $\Delta$ values for each tested combination of discovery algorithm, fitness anprecision metrics using a window size of 100 traces and logs with 2,500 traces. Best results are shadowed in dark and second best results in light grey. In order to compute the $F_{score}$, the error tolerance $\varepsilon$ has been set to a $5\%$ of the log size, that is, half of the trace counbetween real changes.

As can be seen in these results, the parameter with greater impact is the precision metric, being $NP$ precision the one that
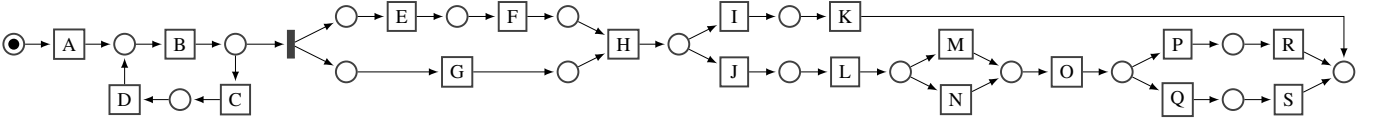
Figure 9. Petri net for the original model [14] used to generate logs. Activity names are shortened for better understandability.

Table II
MEAN $\Delta$ AND $F_{score}$ FOR EVERY TESTED CONFIGURATION

| | | | $F_{score}$ | $\Delta$ |
|---|---|---|---|---|
| IM | $\gamma = AF$ | $\rho = ABA$ | 0.9028 | 4.0915 |
| | | $\rho = NP$ | 0.7985 | 43.6243 |
| | | $\rho = PC$ | 0.9969 | 3.6471 |
| | $\gamma = NR$ | $\rho = ABA$ | 0.9028 | 4.0915 |
| | | $\rho = NP$ | 0.7864 | 41.5318 |
| | | $\rho = PC$ | 0.9969 | 3.5948 |
| | $\gamma = RT$ | $\rho = ABA$ | 0.9425 | 4.0663 |
| | | $\rho = NP$ | 0.7955 | 41.7028 |
| | | $\rho = PC$ | 0.9969 | 3.5948 |
| HM | $\gamma = AF$ | $\rho = ABA$ | 0.9327 | 7.4608 |
| | | $\rho = NP$ | 0.7365 | 36.1014 |
| | | $\rho = PC$ | 0.9789 | 4.4412 |
| | $\gamma = NR$ | $\rho = ABA$ | 0.9402 | 10.3758 |
| | | $\rho = NP$ | 0.7427 | 36.8831 |
| | | $\rho = PC$ | 0.9750 | 5.3828 |
| | $\gamma = RT$ | $\rho = ABA$ | 0.4724 | 6.0812 |
| | | $\rho = NP$ | 0.7025 | 73.5175 |
| | | $\rho = PC$ | 0.7176 | 2.9829 |

| | |
|---|---|
| $IM$ | Inductive Miner |
| $HM$ | Heuristics Miner |
| $AF$ | Alignment Based Fitness |
| $NR$ | Negative Event Recall |
| $RT$ | Percentage of Replayable Traces |
| $ABA$ | Advanced Behavioural Appropriateness |
| $NP$ | Negative Event Precision |
| $PC$ | Precision Change Assessment |



Figure 10. Average execution time for each one of the tested metrics.



Figure 11. Mean $F_{score}$ (higger is better) and $\Delta$ (lower is better) evolution using different window sizes. Shadowed in green is the optimal region.

presents the worst detection resultin terms of both $F_{score}$ and $\Delta$. This metric has a greater sensitivity, so small changes in the frequency of activities in the processed window can dramatically change the precision, and thus bmistaken with structural changes. $ABA$ precision obtains good results in most of cases, but does not finish in some logs, causing some false negatives that affect the final results. Best results are archived using $PC$ precision, as it is insensitive to the activity frequency, counting only if relations appear at least once in the analyzewindow.

Fitness has a smaller impact on results, and all tested metrics obtain very similar results, specially in the case of the $IM$, which ensures a complete fitnessno matter which metric is used to compute it. Only the combination of $RT$ fitness with $HM$ algorithm gets worse results since $HM$ does not ensure a complete trace replayability.

**Impact on computation time.** Figure 10 summarizes the average time required for processing a window of 100 traces by each metric. Test were run on an Intel(R) Core(TM) i7 860 and 16 GiB of RAM. As we can see, $PC$ precision is sensitively faster than the rest of the precision metrics. In the case of fitness, the difference is much lower, being $RT$ the faster one, closely followed by $NR$.

On the basis of the above, the experiments in the following sections have been performed using $IM$ algorithm, $RT$ fitness and $PC$ precision.

*C. Impact of the Window Size*

Figure 11 shows the impact of the window size in both mean $F_{score}$ and mean $\Delta$ for logs of 2,500 traces. The optimal window size for both $F_{score}$ and $\Delta$ ranges between 25 and 100 (this area is boxed in green). Reduced $\Delta$, but with low

$F_{score}$, are obtained for windows with at most 25 traces. In this case, *C2D2* algorithm identifies many false positives because the reference window does not include enough behavior to discover the model. *C2D2* also gets low $F_{score}$ for windows with more than 100 traces because false negatives are identified as the window contains behavior of both before and after the change. In the experiments of the following sections, we will use a window of 100 traces, which is the value that produces the better $F_{score}$ without increasing the $\Delta$.

### D. Comparison with Other Process Drift Detection Algorithms

In this section, *C2D2* algorithm is compared with *Trace-Based ProDrift* (*PD-T*) [14], *Event-Based ProDrift* (*PD-E*) [18] and *TPCDD* [17], the three sudden process drift detection algorithms with best results in the state of the art. Specifically, we used the following configurations:

1) *PD-T* with an *adaptive window* and an initial size of 50;
2) *PD-E* with an *adaptive window* and an initial size of 50. *Relation noise filter threshold* was set to 0 % and *sensitivity* to *very high*, as the authors recommend for analyzing synthetic logs without noise;
3) *TPCDD* with minimum window size set to 100 and *DB-SCAN radius* to 10.

Table III to Table VI show the detailed results. *C2D2* outperforms the rest of algorithms in terms of $F_{score}$ in all the cases, except in OIR log, getting always the best average value. Moreover, *C2D2* is also the second best in terms of delay, very close to *TPCDD*, and being all the values in the same order of magnitude.

For the 2,500 trace logs, *C2D2* obtains a $F_{score}$ of 1.0 in all logs, except in OIR, where it returns a false positive. *TPCDD* also obtains a perfect $F_{score}$, except in lp. In this case the error is due to a false negative (i.e., a change that remains undetected). Both *PD-T* and *PD-E* have much worse results, having 7 and 8 cases where they does not detect all changes. In fact, *PD-T* even is unable to detect any change in cd.

For logs with 5,000, 7,500 and 10,000 traces, *C2D2* gets similar results, and only do not have a perfect $F_{score}$ in OIR, again because of false positives. However, *TPCDD* behaves much worse, being unable to get a perfect $F_{score}$ in any log. The same happens to *PD-E* in logs with 7,500 and 10,000 traces, where it can not detect all changes correctly. *PD-T* gets worse $F_{score}$ in most of the logs, and, in addition, is unable to detect any change in cd and pl, for logs with 5,000 and 7,500 traces, and in cd and cm, for logs with 10,000 traces.

In terms of $\Delta$, *C2D2* and *TPCDD* obtain similar results for all the logs, being able to detect all changes always with less than 10 traces of delay. In comparison with *PD-T* and *PD-E*, *C2D2* $\Delta$ are always an order of magnitude below, thus being able to detect changes closer to the point where they really happened.

As a summary, Figure 12 shows the mean $F_{score}$ and the mean $\Delta$ for each algorithm. As can be seen, *C2D2* outperforms every other algorithm in mean $F_{score}$, while maintaining a very low $\Delta$. Furthermore, $F_{score}$ results have been evaluated using a statistical test in order to confirm the quality of the proposed
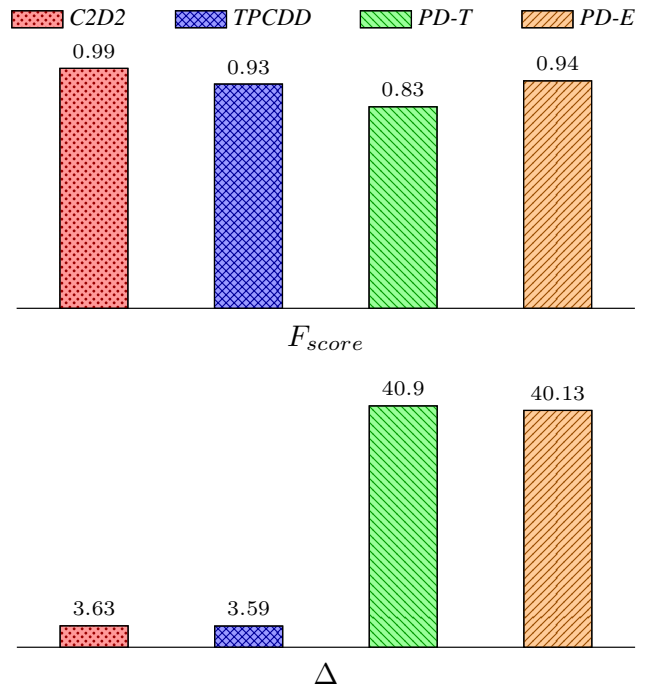


Figure 12. Mean $F_{score}$ and $\Delta$ values for *C2D2*, *PD-T*, *PD-E*, *TPCDD*.

method. A *non-parametric one-vs-all* test has been executed using the STAC tool[2] [26]. Namely, *Friedman's Aligned Ranks* test with a significance level of 0.05 has been used. The test result confirms that $H_0$ is rejected, thus algorithms do not converge to the same mean, being *C2D2* statistically the best in terms of $F_{score}$. Table VII shows the rankings of the algorithms. In this case, *C2D2* clearly outscores the other algorithms.

A post-hoc test (namely *Holm* post-hoc test), a pairwise comparison, was also performed. Results are shown in Table VIII. For all the cases, $H_0$ is rejected, meaning that no algorithm is able to equal *C2D2* in $F_{score}$ values. Finally, tests were also conducted for $\Delta$ values. *C2D2* and *TPCDD* are tied as the best-performing algorithms. *Friedman's Aligned Ranks* for $\Delta$ is depicted in Table IX. Both *C2D2* and *TPCDD* obtain similar scores, thus rejecting the hypothesis that one of the algorithms outperforms the others. Table X shows the post-hoc test for $\Delta$. As we can see in the last row, the null hypothesis $H_0$ between *C2D2* and *TPCDD* cannot be rejected, meaning that both algorithms converge to the same mean, which is consistent with the scores obtained by *Friedman's Aligned Ranks*.

### VI. CONCLUSIONS AND FUTURE WORK

In this paper we presented *C2D2*, an approach to the *offline* detection of *sudden control-flow drifts* in process mining. *C2D2* drift detection is supported by the hypotheses of conformance checking measures being suitable to detect control-flow drifts. Specifically, we argue that fitness and precision are complementary metrics in concept drift, and while fitness is usefull to identify traces that are not supported by the model, precision

---

[2]tec.citius.usc.es/stac

Table III
MEAN $F_{score}$ AND $\Delta$ VALUES FOR EACH ALGORITHM USING LOGS WITH 2,500 TRACES.

| Log | C2D2 | | TPCDD | | PD-T | | PD-E | |
|---|---|---|---|---|---|---|---|---|
| | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ |
| cb | 1.0000 | 6.2222 | 1.0000 | 5.5556 | 0.2000 | 90.0000 | 1.0000 | 78.4444 |
| cd | 1.0000 | 1.7778 | 1.0000 | 0.7778 | 0.0000 | — | 1.0000 | 42.4444 |
| cf | 1.0000 | 3.8889 | 1.0000 | 1.4444 | 1.0000 | 40.0000 | 1.0000 | 49.7778 |
| cm | 1.0000 | 5.8889 | 1.0000 | 11.5556 | 0.7143 | 66.0000 | 0.9412 | 95.1250 |
| cp | 1.0000 | 2.8889 | 1.0000 | 1.2222 | 1.0000 | 45.2222 | 1.0000 | 39.0000 |
| lp | 1.0000 | 3.6667 | 0.9474 | 10.0000 | 1.0000 | 52.7778 | 0.9412 | 36.5000 |
| sw | 1.0000 | 2.8889 | 1.0000 | 1.0000 | 1.0000 | 43.5556 | 0.9412 | 38.1250 |
| pl | 1.0000 | 1.7778 | 1.0000 | 1.3333 | 0.9412 | 36.0000 | 0.9412 | 51.3750 |
| pm | 1.0000 | 3.2222 | 1.0000 | 3.3333 | 0.8000 | 49.6667 | 1.0000 | 48.4444 |
| re | 1.0000 | 2.0000 | 1.0000 | 1.1111 | 1.0000 | 19.8889 | 0.9412 | 49.8750 |
| rp | 1.0000 | 3.0000 | 1.0000 | 1.2222 | 1.0000 | 45.2222 | 1.0000 | 65.0000 |
| IOR | 1.0000 | 2.6667 | 1.0000 | 2.0000 | 1.0000 | 38.7778 | 1.0000 | 63.3333 |
| IRO | 1.0000 | 6.2222 | 1.0000 | 2.7778 | 1.0000 | 52.3333 | 1.0000 | 43.0000 |
| OIR | 0.9474 | 5.7778 | 1.0000 | 0.6667 | 0.6154 | 29.5000 | 0.9412 | 15.7500 |
| ORI | 1.0000 | 2.5556 | 1.0000 | 0.6667 | 1.0000 | 50.5556 | 1.0000 | 35.0000 |
| RIO | 1.0000 | 4.6667 | 1.0000 | 10.8889 | 0.9412 | 51.1250 | 0.9412 | 53.8750 |
| ROI | 1.0000 | 2.0000 | 1.0000 | 0.7778 | 1.0000 | 32.3333 | 0.9412 | 40.1250 |
| AVG | 0.9969 | 3.5948 | 0.9969 | 3.3137 | 0.8360 | 46.4349 | 0.9723 | 49.7173 |

Table IV
MEAN $F_{score}$ AND $\Delta$ VALUES FOR EACH ALGORITHM USING LOGS WITH 5,000 TRACES.

| Log | C2D2 | | TPCDD | | PD-T | | PD-E | |
|---|---|---|---|---|---|---|---|---|
| | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ |
| cb | 1.0000 | 6.7778 | 0.9000 | 6.5556 | 0.7143 | 57.4000 | 1.0000 | 99.6667 |
| cd | 1.0000 | 1.2222 | 0.9474 | 0.7778 | 0.0000 | — | 0.9412 | 2.5555 |
| cf | 1.0000 | 3.3333 | 0.9474 | 9.2222 | 1.0000 | 25.7777 | 1.0000 | 50.1111 |
| cm | 1.0000 | 7.7778 | 0.9474 | 6.1111 | 0.7143 | 51.8000 | 0.8750 | 117.5714 |
| cp | 1.0000 | 3.5556 | 0.9000 | 5.4444 | 1.0000 | 33.3333 | 1.0000 | 66.2222 |
| lp | 1.0000 | 2.0000 | 0.9000 | 4.1111 | 1.0000 | 48.4444 | 1.0000 | 45.3333 |
| sw | 1.0000 | 2.7778 | 0.9474 | 0.6667 | 1.0000 | 29.6667 | 1.0000 | 21.3333 |
| pl | 1.0000 | 1.2222 | 0.9474 | 0.7778 | 0.0000 | — | 1.0000 | 51.2222 |
| pm | 1.0000 | 3.6667 | 0.9474 | 3.4444 | 0.9412 | 37.0000 | 1.0000 | 25.8889 |
| re | 1.0000 | 2.0000 | 0.9474 | 1.5556 | 1.0000 | 19.1111 | 0.9474 | 16.7778 |
| rp | 1.0000 | 2.8889 | 0.9000 | 0.4444 | 1.0000 | 28.2222 | 0.9412 | 48.7500 |
| IOR | 1.0000 | 2.2222 | 0.9474 | 2.7778 | 1.0000 | 24.2222 | 1.0000 | 52.0000 |
| IRO | 1.0000 | 5.4444 | 0.9000 | 1.7778 | 1.0000 | 49.4444 | 0.9412 | 31.2500 |
| OIR | 0.9474 | 7.7778 | 0.9000 | 1.0000 | 1.0000 | 26.6667 | 0.9412 | 0.0000 |
| ORI | 1.0000 | 3.2222 | 0.9000 | 1.1111 | 1.0000 | 36.1111 | 0.9412 | 22.7500 |
| RIO | 1.0000 | 5.2222 | 0.9474 | 1.7778 | 1.0000 | 45.2222 | 1.0000 | 59.7778 |
| ROI | 1.0000 | 2.0000 | 0.9474 | 0.4444 | 1.0000 | 26.0000 | 0.9412 | 7.5000 |
| AVG | 0.9969 | 3.7124 | 0.9279 | 2.8235 | 0.8453 | 35.8948 | 0.9692 | 42.2771 |

looks for behaviour not present in the window of traces. Related to this, we propose the usage of two approximations for each one of the metrics that have a low computational complexity and to detect changes in models.

Our approach has been validated against a synthetic benchmarking dataset formed by 68 logs, outperforming the best concept drift algorithms in terms of accuracy ($F_{score}$) while maintaining a minimum delay ($\Delta$). Finally, we did a statistical test over the results of all the algorithms to confirm that the presented solution is statistically better in terms of accuracy.

As future work, we plan to study the usage of memory mechanisms in order to improve the results. We plan also to extend the method to deal with other types of changes, as the gradual and incremental drifts, and to be executed in *online* environments, where the requirements in terms of computational complexity are different.

Table V

MEAN $F_{score}$ AND $\Delta$ VALUES FOR EACH ALGORITHM USING LOGS WITH 7,500 TRACES.

| Log | C2D2 | | TPCDD | | PD-T | | PD-E | |
|-----|---------|--------|---------|---------|---------|---------|---------|---------|
| | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ |
| **cb** | 1.0000 | 7.7778 | 0.9474 | 16.6667 | 1.0000 | 68.6667 | 0.9474 | 70.2222 |
| **cd** | 1.0000 | 2.1111 | 0.9474 | 1.0000 | 0.0000 | — | 0.9474 | 2.2222 |
| **cf** | 1.0000 | 2.7778 | 0.9474 | 1.5556 | 1.0000 | 22.3333 | 0.9474 | 37.5556 |
| **cm** | 1.0000 | 5.3333 | 0.9474 | 4.4444 | 0.9412 | 83.8750 | 0.8889 | 78.6250 |
| **cp** | 1.0000 | 3.2222 | 0.9000 | 1.3333 | 1.0000 | 33.2222 | 0.9474 | 26.5556 |
| **lp** | 1.0000 | 1.4444 | 0.7500 | 1.2222 | 1.0000 | 52.8889 | 0.7200 | 34.6667 |
| **sw** | 1.0000 | 2.6667 | 0.8571 | 0.8889 | 1.0000 | 32.4444 | 0.9000 | 16.1111 |
| **pl** | 1.0000 | 2.1111 | 0.9474 | 1.6667 | 0.0000 | — | 0.9474 | 30.8889 |
| **pm** | 1.0000 | 3.1111 | 0.9474 | 2.6667 | 1.0000 | 43.4444 | 0.9474 | 14.6667 |
| **re** | 1.0000 | 2.0000 | 0.9474 | 0.5556 | 1.0000 | 21.4444 | 0.9474 | 27.5556 |
| **rp** | 1.0000 | 2.5556 | 0.8571 | 30.6667 | 1.0000 | 29.8889 | 0.9000 | 50.5556 |
| **IOR** | 1.0000 | 4.3333 | 0.9474 | 10.7778 | 0.9000 | 34.5556 | 0.9474 | 53.2222 |
| **IRO** | 1.0000 | 4.2222 | 0.9000 | 2.7778 | 1.0000 | 53.4444 | 0.9474 | 14.3333 |
| **OIR** | 0.7500 | 2.0000 | 0.7826 | 0.3333 | 1.0000 | 69.7778 | 0.9474 | 0.0000 |
| **ORI** | 1.0000 | 2.3333 | 0.8571 | 0.8889 | 1.0000 | 33.1111 | 0.9000 | 36.8889 |
| **RIO** | 1.0000 | 4.2222 | 0.8571 | 2.4444 | 1.0000 | 40.6667 | 0.9000 | 47.6667 |
| **ROI** | 1.0000 | 2.0000 | 0.9474 | 1.0000 | 1.0000 | 31.3333 | 0.9474 | 7.8888 |
| **AVG** | 0.9853 | 3.1895 | 0.8993 | 4.7582 | 0.8730 | 43.4065 | 0.9194 | 32.3309 |

Table VI

MEAN $F_{score}$ AND $\Delta$ VALUES FOR EACH ALGORITHM USING LOGS WITH 10,000 TRACES.

| Log | C2D2 | | TPCDD | | PD-T | | PD-E | |
|-----|---------|--------|---------|---------|---------|---------|---------|---------|
| | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ | $F_{score}$ | $\Delta$ |
| **cb** | 1.0000 | 9.0000 | 0.8571 | 18.2222 | 0.5000 | 41.6667 | 0.8571 | 84.4444 |
| **cd** | 1.0000 | 1.5556 | 0.9474 | 1.0000 | 0.0000 | — | 0.9000 | 3.3333 |
| **cf** | 1.0000 | 3.8889 | 0.9474 | 5.8889 | 1.0000 | 29.6667 | 0.9474 | 45.8889 |
| **cm** | 1.0000 | 8.5556 | 0.9000 | 7.6667 | 0.0000 | — | 0.9474 | 97.7778 |
| **cp** | 1.0000 | 4.4444 | 0.9000 | 2.3333 | 1.0000 | 33.5556 | 0.9000 | 30.7778 |
| **lp** | 1.0000 | 2.3333 | 0.8182 | 8.1111 | 1.0000 | 47.1111 | 0.8182 | 32.8889 |
| **sw** | 1.0000 | 3.6667 | 0.9474 | 1.5556 | 1.0000 | 33.4444 | 0.9000 | 4.6667 |
| **pl** | 1.0000 | 1.5556 | 0.9474 | 1.1111 | 0.2000 | 84.0000 | 0.9474 | 30.7778 |
| **pm** | 1.0000 | 3.6667 | 0.9474 | 2.1111 | 0.7500 | 31.3333 | 0.9474 | 9.0000 |
| **re** | 1.0000 | 2.0000 | 0.9000 | 0.7778 | 1.0000 | 17.7778 | 0.9000 | 23.0000 |
| **rp** | 1.0000 | 3.7778 | 0.8182 | 1.3333 | 1.0000 | 31.2222 | 0.8571 | 46.2222 |
| **IOR** | 1.0000 | 2.7778 | 0.9000 | 3.0000 | 1.0000 | 27.7778 | 0.9474 | 59.5556 |
| **IRO** | 1.0000 | 7.6667 | 0.9474 | 2.1111 | 0.8750 | 50.1429 | 0.8571 | 18.0000 |
| **ORI** | 1.0000 | 3.8889 | 0.8182 | 1.1111 | 1.0000 | 32.6667 | 0.9000 | 54.6666 |
| **OIR** | 0.7826 | 2.0000 | 0.7200 | 0.2222 | 1.0000 | 37.7778 | 0.9474 | 2.7778 |
| **RIO** | 1.0000 | 5.8889 | 0.7500 | 2.1111 | 0.8889 | 41.0000 | 0.8182 | 63.0000 |
| **ROI** | 1.0000 | 2.0000 | 0.9474 | 0.4444 | 1.0000 | 20.6667 | 0.9000 | 8.4444 |
| **AVG** | 0.9872 | 4.0392 | 0.8831 | 3.4771 | 0.7773 | 37.8540 | 0.8995 | 36.1895 |

Table VII

$F_{score}$ ONE-VS-ALL STATISTICAL TEST RESULTS.

| Rank | Algorithm |
|----------|-----------|
| 85.5221 | C2D2 |
| 146.6985 | PD-T |
| 152.2132 | PD-E |
| 162.5662 | TPCDD |

Table VIII

$F_{score}$ POST-HOC STATISTICAL TEST RESULTS.

| Algorithms | $p_{value}$ | Result |
|------------|-------------|--------|
| C2D2 vs. PD-T | $< 1 \times 10^{-10}$ | $H_0$ is rejected |
| C2D2 vs. PD-E | $< 1 \times 10^{-10}$ | $H_0$ is rejected |
| C2D2 vs. TPCDD | $< 1 \times 10^{-10}$ | $H_0$ is rejected |

REFERENCES

[1] W. M. P. van der Aalst *et al.*, "Process Mining Manifesto," in *Business Process Management Workshops - BPM 2011 International Workshops,*

Table IX
Δ ONE-VS-ALL STATISTICAL TEST RESULTS.

| Rank | Algorithm |
|---|---|
| 62.8629 | TPCDD |
| 63.7984 | C2D2 |
| 185.6129 | PD-E |
| 185.7258 | PD-T |

Table X
Δ POST-HOC STATISTICAL TEST RESULTS.

| Algorithms | $p_{value}$ | Result |
|---|---|---|
| C2D2 vs. PD-T | $< 1 \times 10^{-10}$ | $H_0$ is rejected |
| C2D2 vs. PD-E | $< 1 \times 10^{-10}$ | $H_0$ is rejected |
| C2D2 vs. TPCDD | 0.9431 | $H_0$ is accepted |

Clermont-Ferrand, France, August 29, 2011, Revised Selected Papers, Part I, 2011.

[2] J. Gama, I. Zliobaite, A. Bifet, M. Pechenizkiy, and A. Bouchachia, "A survey on concept drift adaptation," *ACM Comput. Surv.*, 2014.

[3] R. P. J. C. Bose, W. M. P. van der Aalst, I. Zliobaite, and M. Pechenizkiy, "Dealing With Concept Drifts in Process Mining," *IEEE Trans. Neural Netw. Learning Syst.*, 2014.

[4] J. Carmona and R. Gavaldà, "Online Techniques for Dealing with Concept Drift in Process Mining," in *Advances in Intelligent Data Analysis XI - 11th International Symposium, IDA 2012, Helsinki, Finland, October 25-27, 2012. Proceedings*, 2012.

[5] P. Weber, B. Bordbar, and P. Tiño, "Real-Time Detection of Process Change using Process Mining," in *2011 Imperial College Computing Student Workshop, ICCSW 2011, London, United Kingdom, September 29-30, 2011. Proceedings.*, 2011.

[6] S. B. Junior, G. M. Tavares, V. G. T. da Costa, P. Ceravolo, and E. Damiani, "A Framework for Human-in-the-loop Monitoring of Concept-drift Detection in Event Log Stream," in *Companion of the The Web Conference 2018 on The Web Conference 2018, WWW 2018, Lyon , France, April 23-27, 2018*, 2018.

[7] F. M. Maggi, A. Burattin, M. Cimitile, and A. Sperduti, "Online Process Discovery to Detect Concept Drifts in LTL-Based Declarative Process Models," in *On the Move to Meaningful Internet Systems: OTM 2013 Conferences - Confederated International Conferences: CoopIS, DOA-Trusted Cloud, and ODBASE 2013, Graz, Austria, September 9-13, 2013. Proceedings*, 2013.

[8] J. Martjushev, R. P. J. C. Bose, and W. M. P. van der Aalst, "Change Point Detection and Dealing with Gradual and Multi-order Dynamics in Process Mining," in *Perspectives in Business Informatics Research - 14th International Conference, BIR 2015, Tartu, Estonia, August 26-28, 2015, Proceedings*, 2015.

[9] M. K. M. V., L. Thomas, and A. Basava, "Capturing the Sudden Concept Drift in Process Mining," in *Proceedings of the International Workshop on Algorithms & Theories for the Analysis of Event Data, ATAED 2015, Satellite event of the conferences: 36th International Conference on Application and Theory of Petri Nets and Concurrency Petri Nets 2015 and 15th International Conference on Application of Concurrency to System Design ACSD 2015, Brussels, Belgium, June 22-23, 2015.*, 2015.

[10] T. Li, T. He, Z. Wang, Y. Zhang, and D. Chu, "Unraveling Process Evolution by Handling Concept Drifts in Process Mining," in *2017 IEEE International Conference on Services Computing, SCC 2017, Honolulu, HI, USA, June 25-30, 2017*, 2017.

[11] R. Accorsi and T. Stocker, "Discovering Workflow Changes with Time-Based Trace Clustering," in *Data-Driven Process Discovery and Analysis - First International Symposium, SIMPDA 2011, Campione d'Italia, Italy, June 29 - July 1, 2011, Revised Selected Papers*, 2011.

[12] D. Luengo and M. Sepúlveda, "Applying Clustering in Process Mining to Find Different Versions of a Business Process That Changes over Time," in *Business Process Management Workshops - BPM 2011 International Workshops, Clermont-Ferrand, France, August 29, 2011, Revised Selected Papers, Part I*, 2011.

[13] B. Hompes, J. C. A. M. Buijs, W. M. P. van der Aalst, P. Dixit, and H. Buurman, "Detecting Change in Processes Using Comparative Trace Clustering," in *Proceedings of the 5th International Symposium on Data-driven Process Discovery and Analysis (SIMPDA 2015), Vienna, Austria, December 9-11, 2015.*, 2015.

[14] A. Maaradji, M. Dumas, M. L. Rosa, and A. Ostovar, "Detecting Sudden and Gradual Drifts in Business Processes from Execution Traces," *IEEE Trans. Knowl. Data Eng.*, 2017.

[15] G. T. Lakshmanan, P. T. Keyser, and S. Duan, "Detecting changes in a semi-structured business process through spectral graph analysis," in *Workshops Proceedings of the 27th International Conference on Data Engineering, ICDE 2011, April 11-16, 2011, Hannover, Germany*, 2011.

[16] A. Seeliger, T. Nolle, and M. Mühlhäuser, "Detecting Concept Drift in Processes using Graph Metrics on Process Graphs," in *Proceedings of the 9th Conference on Subject-oriented Business Process Management, S-BPM ONE 2017, Darmstadt, Germany, March 30-31, 2017*, 2017.

[17] C. Zheng, L. Wen, and J. Wang, "Detecting Process Concept Drifts from Event Logs," in *On the Move to Meaningful Internet Systems. OTM 2017 Conferences - Confederated International Conferences: CoopIS, C&TC, and ODBASE 2017, Rhodes, Greece, October 23-27, 2017, Proceedings, Part I*, 2017.

[18] A. Ostovar, A. Maaradji, M. L. Rosa, A. H. M. ter Hofstede, and B. F. van Dongen, "Detecting drift from event streams of unpredictable business processes," in *ER*, 2016.

[19] B. Weber, M. Reichert, and S. Rinderle-Ma, "Change patterns and change support features - Enhancing flexibility in process-aware information systems," *Data Knowl. Eng.*, 2008.

[20] A. Bifet and R. Gavaldà, "Learning from Time-Changing Data with Adaptive Windowing," in *Proceedings of the Seventh SIAM International Conference on Data Mining, April 26-28, 2007, Minneapolis, Minnesota, USA*, 2007.

[21] R. P. J. C. Bose and W. M. P. van der Aalst, "Trace Clustering Based on Conserved Patterns: Towards Achieving Better Process Models," in *Business Process Management Workshops, BPM 2009 International Workshops, Ulm, Germany, September 7, 2009. Revised Papers*, 2009.

[22] A. Rozinat and W. M. P. van der Aalst, "Conformance checking of processes based on monitoring real behavior," *Inf. Syst.*, 2008.

[23] S. K. L. M. vanden Broucke, J. D. Weerdt, J. Vanthienen, and B. Baesens, "Determining process model precision and generalization with weighted artificial negative events," *IEEE Trans. Knowl. Data Eng.*, 2014.

[24] S. J. J. Leemans, D. Fahland, and W. M. P. van der Aalst, "Discovering Block-Structured Process Models from Event Logs - A Constructive Approach," in *Application and Theory of Petri Nets and Concurrency - 34th International Conference, PETRI NETS 2013, Milan, Italy, June 24-28, 2013. Proceedings*, 2013.

[25] A. J. M. M. Weijters and J. T. S. Ribeiro, "Flexible Heuristics Miner (FHM)," in *Proceedings of the IEEE Symposium on Computational Intelligence and Data Mining, CIDM 2011, part of the IEEE Symposium Series on Computational Intelligence 2011, April 11-15, 2011, Paris, France*, 2011.

[26] I. Rodríguez-Fdez, A. Canosa, M. Mucientes, and A. Bugarín, "STAC: A web platform for the comparison of algorithms using statistical tests," in *FUZZ-IEEE*, 2015.