

Learning what to Monitor: using Machine Learning to Improve Past STL Monitoring

Andrea Brunello, Luca Geatti, Angelo Montanari and Nicola Saccomanno

University of Udine, Udine, Italy

{andrea.brunello, luca.geatti, angelo.montanari, nicola.saccomanno}@uniud.it

Abstract

Monitoring is a runtime verification technique that can be used to check whether an execution of a system (*trace*) satisfies or not a given set of properties. Compared to other formal verification techniques, e.g., model checking, one needs to specify the properties to be monitored, but a complete model of the system is no longer necessary. First, we introduce the *pure past* fragment of *Signal Temporal Logic* (ppSTL), and we use it to define the monitorable safety (G(ppSTL)) and cosafety (F(ppSTL)) fragments of STL, which properly extend the commonly-used bounded-future fragment. Then, we devise a multi-objective genetic programming algorithm to automatically extend the set of properties to monitor on the basis of the history of failure traces collected over time. The framework resulting from the integration of the monitor and the learning algorithm is then experimentally validated on various public datasets. The outcomes of the experimentation confirm the effectiveness of the proposed solution.

1 Introduction

In this paper, we propose a tight integration of monitoring and machine learning for preemptive failure detection, with formal guarantees on interpretability, monitorability, and expressiveness.

Monitoring is a runtime verification technique for the analysis of complex systems [Leucker and Schallhart, 2009]. It consists of the generation of a *monitor* that is paired with the system under analysis, and it reports a positive (resp., negative) verdict whenever the current execution, and all its continuations, are guaranteed to be good (resp., bad). Therefore, verdicts of monitors are always *irrevocable*. Expressing good and bad behaviors to be monitored is done by means of temporal logics, in particular *Signal Temporal Logic* (STL) which proved to be quite effective in this context [Maler and Nickovic, 2004]. Not all formulas of STL are *monitorable*: there are formulas for which it is impossible to produce a verdict

by a finite number of observations [Bauer *et al.*, 2011]. Although not being able to check exhaustively all possible executions of a system, like, e.g., model checking, monitoring offers a number of advantages, including: (i) it can be applied directly on the implementation of the system, avoiding the risk of modeling errors; (ii) monitoring algorithms are usually very fast. However, an important limitation, which severely affects the effectiveness of monitoring, remains. Modern systems possess such a level of complexity that it is impossible for a system engineer to specify in advance all properties to be monitored. Moreover, even when this is possible, minor changes to the system may introduce unforeseen bugs.

In this paper, we investigate how to solve this problem by pairing monitoring with machine learning, which is used to learn in an iterative fashion new formulas to monitor by analysing trace prefixes that lead to failure. The method that we propose has three distinguishing features:

- *interpretability*: the machine learning methods that we use manipulate and produce only STL formulae, that can be easily inspected by a system engineer;
- formal guarantees on *monitorability*: every learned formula is guaranteed to be monitorable;
- *expressiveness*: the language for the specification of properties is proved to be able to express more properties than languages typically used in this context.

Our contributions are the following. First, we focus on the specification language to be used to specify properties. We introduce G(ppSTL) and F(ppSTL), the syntactic *safety* and *cosafety* fragments of STL, respectively. A formula of G(ppSTL) is of the form $G(\alpha)$ where G is the *globally* operator (which forces α to be always true) and α is a formula looking only to the past. Therefore, formulas of G(ppSTL) are used to express conditions that has to hold always, *i.e.*, invariants. Similarly, F(ppSTL) is used to express properties that should hold at least one time in the future, like, e.g., planning goals. A key feature of the two fragments is that they can look arbitrarily back into the past, a feature which is not offered by the fragments of STL that are commonly used in this context, in particular the *bounded future fragment* of STL (bfSTL). In the following, we formally prove that these fragments are *more expressive* than bfSTL.

We also give formal guarantees on the *monitorability* of the formulas of the fragments: we prove that each formula

Appendix and Supplementary materials are available at: <https://github.com/dslab-uniud/ppSTL-IJCAI2024>

in $G(\text{ppSTL})$ and $F(\text{ppSTL})$ is monitorable. On the one hand, this avoids the risk of generating nonmonitorable formulas in the learning phase that we will describe later; on the other hand, checking monitorability, which can be cumbersome [Havelund and Peled, 2023], is no more necessary. In addition, for the case of formulas interpreted over qualitative time, *i.e.*, system's executions with no timestamps attached to each time point, we prove that $G(\text{ppSTL})$ and $F(\text{ppSTL})$ capture the *entire* class of safety and cosafety properties, which form an important subclass of monitorable properties.

Third, we devise a framework for the automatic discovery of relevant properties, written in $G(\text{ppSTL})$ and $F(\text{ppSTL})$, based on traces that lead to failure. The objective here is to generate a pool of formulas that helps *anticipating* the identification of failures, a task carried out by a multi-objective genetic algorithm. The results of our experiments show that the pool of formulas obtained in this way can be effectively used for anticipating the detection of failures. Moreover, they can be easily checked by a system engineer, making the *interpretability* of results a distinguished feature of our method. In a dedicated section, we discuss the differences of our methodology with respect to existing solutions for the integration of monitoring and machine learning.

The paper is structured as follows. In Section 2, we provide some background knowledge. Section 3 is dedicated to the safety and cosafety fragments of STL and to their theoretical properties (monitorability and expressiveness). In Section 4, we describe our methodology to learn new formulas based on genetic programming. The outcomes of the experimental evaluation are reported in Section 5. In Section 6, we discuss the work done, highlight its strength and some limitations still present, and provide future research directions.

2 Background

2.1 Signal and Metric Temporal Logic

Let \mathbb{T} be a set of *timestamps*, which are numbers attached to each time point that represent the (real) time at which an event has occurred. There are several possible choices for \mathbb{T} , depending on how time instants are modeled, mainly: (i) $\mathbb{T} := \mathbb{N}$, for *qualitative-time*; (ii) $\mathbb{T} := \mathbb{R}$, for *real-time*.

Syntax

From now on, given a set of variables x_1, \dots, x_n , we denote by \mathbb{D} their domain. We define the syntax of *Signal Temporal Logic* [Maler and Nickovic, 2004] as follows.

Definition 1 (Signal Temporal Logic). *Formulas ϕ of Signal Temporal Logic (STL, for short) are inductively defined as follows:*

$$\phi := f_i(\bar{x}) \otimes c \mid \neg\phi \mid \phi \vee \phi \mid X\phi \mid \phi U_I \phi \mid Y\phi \mid \phi S_I \phi$$

where $\bar{x} := \langle x_1, \dots, x_n \rangle$ for some $n \in \mathbb{N}$, each variable x_j , for $1 \leq j \leq n$, takes value in the set \mathbb{D} , $f_i : \mathbb{D}^n \rightarrow \mathbb{D}$ is a computable function, $c \in \mathbb{D}$, and $\otimes \subseteq \mathbb{D} \times \mathbb{D}$. In addition, we require I to be an interval of the form $[a, b]$, with $a \leq b$, or $[a, \infty)$, where $a, b \in \mathbb{T}$ are two timestamps represented in binary notation.

Formulas of type $f_i(\bar{x}) \otimes c$ are called *atomic formulas*. Modalities U_I and S_I are called *until* and

since, respectively. We define the standard shortcut operators as follows: (i) *true*: $\top := f_i(\bar{x}) \vee \neg f_i(\bar{x})$; (ii) *weak tomorrow*: $X\phi := \neg X\neg\phi$; (iii) *eventually*: $F_I\phi := \top U_I \phi$; (iv) *globally*: $G_I\phi := \neg F_I\neg\phi$; (v) *release*: $\phi_1 R_I \phi_2 := \neg(\neg\phi_1 U_I \neg\phi_2)$; (vi) *weak yesterday*: $\bar{Y}\phi := \neg Y\neg\phi$; (vii) *once*: $O_I\phi := \top S_I \phi$; (viii) *historically*: $H_I\phi := \neg O_I\neg\phi$; (ix) *triggers*: $\phi_1 T_I \phi_2 := \neg(\neg\phi_1 S_I \neg\phi_2)$.

When $I = [0, \infty)$, we say that U_I (resp., S_I) is *unbounded* and we simply write U (resp., S); otherwise, it is *bounded*. The same holds for the other operators. Modalities X , U_I , and all shortcuts derived from them, are called *future modalities*, while Y , S_I , and all shortcuts derived from them, are called *past modalities*. We denote by STL the set of STL formulas.

The *bounded future fragment* of STL (denoted by bfSTL) is the set of STL formulas such that each of its temporal modalities is bounded by an interval of the form $[a, b]$, with $a, b \in \mathbb{T}$.

Metric Temporal Logic (denoted by MTL) is the set of formulas obtained from the same grammar as STL, but by replacing the base case, *i.e.*, $f_i(\bar{x}) \otimes c$, by an atomic proposition p taken from a set \mathcal{AP} .

Semantics

From now on, we fix a set of n variables x_1, \dots, x_n . A *real-time n -valued state* is a pair (t, v) , where $t \in \mathbb{T}$ is a timestamp and $v \in \mathbb{D}^n$ represents the n values d_1, \dots, d_n for the variables x_1, \dots, x_n , respectively. A *real-time n -valued trace* σ is a sequence of real-time n -valued states: σ is *infinite* when $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^\omega$ and *finite* when $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^*$. We often write σ as the sequence $\langle \sigma_0, \sigma_1, \dots \rangle$. We denote by $|\sigma|$ the *length* of σ , *i.e.*, the number of its states; if σ is infinite, then we set $|\sigma| = \omega$. For any $0 \leq i < |\sigma|$, we denote by $\sigma_{[0, i]}$ the *prefix* of σ up to position i . For all $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^*$ and all $\sigma' \in (\mathbb{T} \times \mathbb{D}^n)^* \cup (\mathbb{T} \times \mathbb{D}^n)^\omega$, we denote by $\sigma \cdot \sigma'$ the trace obtained by concatenating σ' to σ .

We interpret STL formulas with variables x_1, \dots, x_n over *infinite real-time n -valued traces* $\sigma = \langle \sigma_0, \sigma_1, \dots, \sigma_i, \dots \rangle$ that have to satisfy the following conditions:

- *(strict) monotonicity*: for all $i, j \in \mathbb{N}$, with $i < j$, if $\sigma_i = (t, v)$ and $\sigma_j = (t', v')$, then $t < t'$;
- *progress*: for all $n \in \mathbb{N}$, there exists $i \in \mathbb{N}$ such that $\sigma_i = (t, v)$ and $n < t$.

A real-time n -valued language \mathcal{L} is a set of real-time n -valued traces. By $\bar{\mathcal{L}}$ we denote the *complement* of \mathcal{L} , that is, $\bar{\mathcal{L}} := \{\sigma \in (\mathbb{T} \times \mathbb{D}^n)^\omega \mid \sigma \notin \mathcal{L}\}$. We say that \mathcal{L} is a *real-time multi-valued language* if and only if it is a real-time n -valued language, for some $n \in \mathbb{N} \setminus \{0\}$.

Given an STL formula ϕ with variables x_1, \dots, x_n , an *infinite real-time n -valued trace* $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^\omega$, and a position i , with $i \geq 0$, we define the *satisfaction* of ϕ over σ at position i , written $\sigma, i \models \phi$, inductively as follows:

1. $\sigma, i \models f_i(\bar{x}) \otimes c$ iff $\sigma_i = (t, (v_1, \dots, v_n))$ and $f_i(v_1, \dots, v_n) \otimes c$ is true;
2. $\sigma, i \models \neg\phi$ iff $\sigma, i \not\models \phi$;
3. $\sigma, i \models \phi_1 \vee \phi_2$ iff $\sigma, i \models \phi_1$ or $\sigma, i \models \phi_2$;
4. $\sigma, i \models X\phi_1$ iff $\sigma, i+1 \models \phi_1$;
5. $\sigma, i \models \phi_1 U_{[a, b]} \phi_2$ iff there exists a $j \geq i$ such that:
 - (i) $\sigma, j \models \phi_2$ and $t + a \leq t' \leq t + b$, where $\sigma_i = (t, v)$

and $\sigma_j = (t', v')$, and (ii) for all $i \leq k < j$, it holds that $\sigma, k \models \phi_1$;

6. $\sigma, i \models Y\phi_1$ iff $i > 0$ and $\sigma, i-1 \models \phi$;
7. $\sigma, i \models \phi_1 S_{[a,b]} \phi_2$ iff there exists a $0 \leq j \leq i$ such that:
 - (i) $\sigma, j \models \phi_2$ and $t-b \leq t' \leq t-a$, where $\sigma_i = (t, v)$ and $\sigma_j = (t', v')$, and (ii) for all $j < k \leq i$, it holds that $\sigma, k \models \phi_1$;

Formulas of STL are interpreted at the *first* position of the trace σ : we say that σ is a *model* of the STL formula ϕ (written $\sigma \models \phi$) iff $\sigma, 0 \models \phi$. We define the language of an STL formula as follows.

Definition 2 (Language of an STL formula). *Given an STL formula ϕ with variables x_1, \dots, x_n , the language of ϕ , denoted by $\mathcal{L}(\phi)$, is the real-time n -valued language $\{\sigma \in (\mathbb{T} \times \mathbb{D}^n)^\omega \mid \sigma \models \phi\}$.*

The semantics of MTL is defined by considering infinite real-time traces in $(\mathbb{T} \times 2^{\mathcal{AP}})^\omega$, for a given set of atomic propositions \mathcal{AP} , instead of real-time multi-valued traces, like for STL. The base case of the semantics is defined as follows:

1. $\sigma, i \models p$ iff $\sigma_i = (t, v)$ and $p \in v$.

The rest of the semantic clauses remain the same as in the case of STL. The language of an MTL formula ϕ over the propositional atoms \mathcal{AP} , denoted by $\mathcal{L}(\phi)$, is defined as $\mathcal{L}(\phi) := \{\sigma \in (\mathbb{T} \times 2^{\mathcal{AP}})^\omega \mid \sigma, 0 \models \phi\}$.

2.2 Safety and cosafety fragments

Cosafety real-time multi-valued languages are defined as follows.

Definition 3 (Cosafety real-time multi-valued language). *Let $\mathcal{L} \subseteq (\mathbb{T} \times \mathbb{D}^n)^\omega$ be a real-time multi-valued language. We say that \mathcal{L} is a cosafety language iff, for all $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^\omega$, if $\sigma \in \mathcal{L}$, then there exists $i \geq 0$ such that $\sigma_{[0,i]} \cdot \sigma' \in \mathcal{L}$, for all $\sigma' \in (\mathbb{T} \times \mathbb{D}^n)^\omega$.*

Safety real-time multi-valued languages are defined as the duals of cosafety ones.

Definition 4 (Safety real-time multi-valued language). *Let $\mathcal{L} \subseteq (\mathbb{T} \times \mathbb{D}^n)^\omega$ be a real-time multi-valued language. We say that \mathcal{L} is a safety language iff $\bar{\mathcal{L}}$ is a cosafety real-time multi-valued language.*

2.3 Monitoring

Monitoring is a lightweight runtime verification technique [Leucker and Schallhart, 2009] and it consists of the generation of a *monitor* that checks an execution of a system either in an online fashion (*i.e.*, at runtime) or in an offline fashion (*e.g.*, by analysing the log of the system). A monitor reports two types of results: an inconclusive output (denoted by $?$) or an irrevocable verdict, in particular a violation (*resp.*, a satisfaction) in case all the continuations of that execution are bad (*resp.*, good). We define a *monitor for a language* $\mathcal{L} \subseteq (\mathbb{T} \times \mathbb{D}^n)^\omega$ as a function $\text{mon}_{\mathcal{L}} : (\mathbb{T} \times \mathbb{D}^n)^* \rightarrow \{\top, \perp, ?\}$ such that, for all $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^*$,

$$\text{mon}_{\mathcal{L}}(\sigma) := \begin{cases} \top & \text{iff } \forall \sigma' \in (\mathbb{T} \times \mathbb{D}^n)^\omega . \sigma \cdot \sigma' \in \mathcal{L} \\ \perp & \text{iff } \forall \sigma' \in (\mathbb{T} \times \mathbb{D}^n)^\omega . \sigma \cdot \sigma' \notin \mathcal{L} \\ ? & \text{otherwise} \end{cases}$$

Given an STL formula ϕ , we will denote with mon_{ϕ} the monitor $\text{mon}_{\mathcal{L}(\phi)}$. When the monitor returns \top , we know that all continuations are good w.r.t. the property ϕ : this is the case, *e.g.*, of planning problems [Ghallab *et al.*, 2004], where ϕ expresses the achievement of a goal. If the monitor returns \perp , we know that there has been an irremediable violation of ϕ : this is the case, for instance, of invariance properties, requiring that something bad never happens [Kupferman and Vardi, 2001].

We point out the similarity between the \top (*resp.*, \perp) result of monitors and cosafety (*resp.*, safety) properties. In fact, all cosafety properties and all safety properties are *monitorable*, in the sense that, for every trace σ , there exists at least one continuation σ' such that $\text{mon}_{\phi}(\sigma \cdot \sigma') \in \{\top, \perp\}$. We define *monitorability* as follows.

Definition 5 (Monitorability). *For each $\mathcal{L} \subseteq (\mathbb{T} \times \mathbb{D}^n)^\omega$, we say that \mathcal{L} is monitorable iff, for all $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^*$, there exists a $\sigma' \in (\mathbb{T} \times \mathbb{D}^n)^*$, such that $\text{mon}_{\phi}(\sigma \cdot \sigma') \neq ?$.*

We say that an STL formula ϕ is monitorable iff $\mathcal{L}(\phi)$ is monitorable. Not all STL formulas are monitorable. As an example, the formula $G(x > 0 \rightarrow Fy < 0)$ stating that every time x is greater than 0 there exists a point in the future where y is negative, is *not* monitorable. However, every (co)safety property is monitorable.

Proposition 1 ([Bauer *et al.*, 2011]). *For every $\mathcal{L} \subseteq (\mathbb{T} \times \mathbb{D}^n)^\omega$, if \mathcal{L} is safety or cosafety, then \mathcal{L} is monitorable.*

We point out that the *vice versa* of Proposition 1 does not hold: there exist monitorable languages that are neither safety nor cosafety [Bauer *et al.*, 2011].

3 Monitorable fragments of STL

In this section, we define the safety and the cosafety syntactic fragments of STL. We show that they express, respectively, only safety and cosafety languages, and thus only *monitorable* properties. In addition, in the case of *qualitative time*, we prove also the *vice versa*, *i.e.*, all safety and all cosafety languages definable in STL can be defined in (one of) the two fragments. The proof of all lemmas and theorems can be found in the Appendix A.

3.1 The G(ppSTL) and the F(ppSTL) fragments

The safety and the cosafety syntactic fragments of STL are based on the *pure past fragment* of STL, which comprises formulas of STL that can look only into the past, and is defined as follows.

Definition 6 (The pure past fragment of STL). *The pure past fragment of STL, denoted by ppSTL, is the set of STL formulas devoid of future operators.*

Unlike the case of STL formulas, in the case of ppSTL we consider only *finite, nonempty* real-time multi-valued traces. Formulas of ppSTL are interpreted at the *last* position of a trace $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^+$, and we say that σ is a *model* of the ppSTL formula ϕ iff $\sigma, |\sigma| - 1 \models \phi$.

We define the syntactic safety and cosafety fragments of STL as follows.

Definition 7 (Safety and Cosafety syntactic fragments of STL). We define the *safety* (resp., *cosafety*) syntactic fragment of STL, denoted by $G(\text{ppSTL})$ (resp., $F(\text{ppSTL})$), as the set of STL formulas of the form $G(\psi)$ (resp., $F(\psi)$), where ψ is a formula of ppSTL.

There are three distinguishing features of the syntax of $G(\text{ppSTL})$ and $F(\text{ppSTL})$ that are worth pointing out:

1. The use of *past modalities*, which allows one to avoid the risk of considering *non-monitorable* formulas.
2. The use of *unbounded intervals*, which allows formulas of these fragments to constrain arbitrarily long (yet finite) traces, in contrast to formulas of the bounded fragment of STL which are able to constrain only finite and *bounded* portions of a trace.
3. The use multi-variable functions. In [Brunello *et al.*, 2023], only functions with arity 1 were allowed (e.g., $x + 1 \leq 3$). Here, we deal with multi-variable functions to allow, for instance, the specification of constraints of the form $|x_1 - x_2| > 0$. This is in line with [Maler and Nickovic, 2004; Donzé and Maler, 2010].

Examples

Here, we give some examples of $G(\text{ppSTL})$ and $F(\text{ppSTL})$ formulas. They concern an *arbiter* that has to give some resources (grants) to the processes that make a request. We suppose to have two variables x_g and x_r , with domain $\mathbb{D} := \mathbb{N}$, that are set to a value strictly greater than 0 iff the arbiter gives a grant and the processes perform a request, respectively.

The simple requirement that “a grant is always preceded by a request in at least 2.1 and at most 7.4 time units” is captured by the formula $\phi := G(x_g > 0 \rightarrow O_{[2.1, 7.4]}(x_r > 0))$. Suppose we want to express the *unbounded* version of the previous requirement, that is, “a grant is always preceded by a request”. The $G(\text{ppSTL})$ formula for this requirement is $G(x_g > 0 \rightarrow YO(x_r > 0))$.

The requirement “there is at least one request followed by a grant” is modelled by the $F(\text{ppSTL})$ formula $F(x_g > 0 \wedge YOx_r > 0)$. Other examples are reported in the Appendix B.

Comparison with bfSTL

In the following, we prove that $G(\text{ppSTL})$ and $F(\text{ppSTL})$ are more expressive than bfSTL (the bounded fragment). The rationale is based on the fact that, while bfSTL formulas can constrain only bounded intervals of a trace, $G(\text{ppSTL})$ and $F(\text{ppSTL})$ can do the same but also over intervals of unbounded length. To show it, we first prove that ppSTL is more expressive than bfSTL.

Proposition 2 (ppSTL is more expressive than bfSTL). *There exists a language $\mathcal{L} \subseteq (\mathbb{T} \times \mathbb{D}^n)^*$ (for some $n \in \mathbb{N}$) such that:*

- *there exists a ppSTL formula ϕ such that $\mathcal{L}(\phi) = \mathcal{L}$;*
- *there exists no bfSTL formula ψ such that $\mathcal{L}(\psi) = \mathcal{L}$.*

We define $G(\text{bfSTL})$ as the set of formulas of the form $G(\alpha)$, where $\alpha \in \text{bfSTL}$, and we define $F(\text{bfSTL})$ analogously. From the previous proposition, and from the fact that $G(\text{bfSTL})$ and $F(\text{bfSTL})$ are syntactic fragments of $G(\text{ppSTL})$ and $F(\text{ppSTL})$, respectively, it follows that:

- $G(\text{ppSTL})$ is strictly more expressive than $G(\text{bfSTL})$;
- $F(\text{ppSTL})$ is strictly more expressive than $F(\text{bfSTL})$.

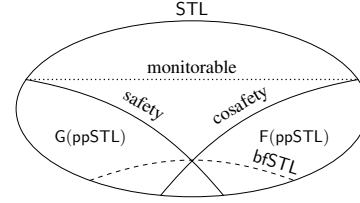


Figure 1: Summary of $G(\text{ppSTL})$ and $F(\text{ppSTL})$ expressive power for qualitative time.

Monitorability of $G(\text{ppSTL})$ and $F(\text{ppSTL})$

In the following, we show that $G(\text{ppSTL})$ (resp., $F(\text{ppSTL})$) expresses only safety (resp., cosafety) real-time multi-valued languages (cf. Definitions 3 and 4).

Lemma 1. *For all $\phi \in G(\text{ppSTL})$ (resp., $\phi \in F(\text{ppSTL})$), it holds that $\mathcal{L}(\phi)$ is a safety (resp., cosafety) real-time multi-valued language.*

It follows from Proposition 1 that all properties definable in $G(\text{ppSTL})$ or in $F(\text{ppSTL})$ are *monitorable*.

Theorem 1. *For all $\phi \in G(\text{ppSTL}) \cup F(\text{ppSTL})$, it holds that $\mathcal{L}(\phi)$ is monitorable.*

In addition, we prove that, when interpreted over *qualitative time*, the fragments $G(\text{ppSTL})$ and $F(\text{ppSTL})$ are *expressively complete*, that is, every safety (resp., cosafety) multi-valued language that can be defined by an STL formula is definable also in $G(\text{ppSTL})$ (resp., in $F(\text{ppSTL})$).

Theorem 2 (Expressive Completeness over qualitative-time). *For all multi-valued languages $\mathcal{L} \subseteq (\mathbb{N} \times \mathbb{D}^n)^*$ definable in STL, it holds that:*

- \mathcal{L} is safety iff there exists a formula ϕ of $G(\text{ppSTL})$ such that $\mathcal{L}(\phi) = \mathcal{L}$;
- \mathcal{L} is cosafety iff there exists a formula ϕ of $F(\text{ppSTL})$ such that $\mathcal{L}(\phi) = \mathcal{L}$.

We point out that, to the best of our knowledge, it is unknown whether Theorem 2 holds for the case of real-time. A recap of the expressiveness results is given in Fig. 1.

4 Failure detection framework

In this section, we describe the proposed framework and the formula learning algorithm that it uses.

4.1 The overall framework

Algorithm 1 describes the framework training phase. It monitors, one after the other, all available failure system traces and, for each one, it simulates its point-by-point arrival (i.e., all the prefixes of the trace). At the end, it returns a pool of formulas \mathcal{P} characterizing bad behaviours of the system.

The procedure gets, as its input, a pool \mathcal{P} of formulas encoding bad behaviours. The pool may be empty, or it may already include some formulas, if they were previously defined by domain experts. In addition, a training set \mathcal{X} is provided, consisting of pairs $(\sigma, is_failure)$, where σ represents a system execution trace and *is_failure* its corresponding label

Algorithm 1: Framework training phase

input: \mathcal{P} initial (possibly empty) pool of formulas,
 \mathcal{X} dataset of labelled traces $(\sigma, is_failure)$,
 q quality requirements

- 1: $\Sigma_{\top} \leftarrow \{\sigma \mid (\sigma, is_failure) \in \mathcal{X} \wedge is_failure = \top\}$
- 2: $\Sigma_{\perp} \leftarrow \{\sigma \mid (\sigma, is_failure) \in \mathcal{X} \wedge is_failure = \perp\}$
- 3: $\Sigma_{\perp} \leftarrow \text{GENAUGMENTEDTRACES}(\Sigma_{\perp})$
- 4: **for** $\sigma \in \Sigma_{\top}$ **do**
- 5: $failure_detected \leftarrow \perp$
- 6: **for** $i \leftarrow 0$ **to** $|\sigma| - 1$ **do**
- 7: $\mathcal{F} \leftarrow \{\psi \in \mathcal{P} \mid \text{mon}_{\psi}(\sigma_{[0,i]}) = \perp\}$
- 8: **if** $\mathcal{F} \neq \emptyset$ **then**
- 9: $failure_detected \leftarrow \top$
- 10: $\mathcal{T} \leftarrow \text{GENAUGMENTEDTRACES}(\sigma_{[0,i]})$
- 11: $\Phi \leftarrow \text{LEARNFORMULA}(\mathcal{T}, \Sigma_{\perp}, q)$
- 12: $\mathcal{P} \leftarrow \mathcal{P} \cup \Phi$
- 13: **break**
- 14: **end if**
- 15: **end for**
- 16: **if** $failure_detected = \perp$ **then**
- 17: $\mathcal{T} \leftarrow \text{GENAUGMENTEDTRACES}(\sigma)$
- 18: $\Phi \leftarrow \text{LEARNFORMULA}(\mathcal{T}, \Sigma_{\perp}, q)$
- 19: $\mathcal{P} \leftarrow \mathcal{P} \cup \Phi$
- 20: **end if**
- 21: **end for**
- 22: **return** \mathcal{P}

(\top , if σ is a trace ending with a failure; \perp otherwise).¹ Finally, quality requirements q are considered that should be satisfied by any new formula extracted during the training process. They will be described more in detail later.

On Line 1 (resp., Line 2) the subset of failure (resp., good) traces is extracted from \mathcal{X} . For each good trace, n_{aug} variants (a framework global parameter) are generated by adding random Gaussian noise as a counter-overfitting measure. As the result, a set of $|\Sigma_{\perp}| + n_{aug} * |\Sigma_{\perp}|$ traces is obtained.

At this point, the framework starts its iterative part, during which a failure trace σ is monitored sequentially, point-by-point (Lines 6–15). At each iteration, the framework restricts its attention to the prefix $\sigma_{[0,i]}$ of trace σ , and it computes the set \mathcal{F} of formulas leading to a *violation* (Line 7). To such an extent, it executes the monitoring tool `rtamt` [Ničković and Yamaguchi, 2020]. Since all formulas $\psi \in \mathcal{P}$ are meant to encode bad behaviors, we say that a formula ψ leads to a violation if $\text{mon}_{\psi}(\sigma_{[0,i]}) = \perp$ (`mon` is defined as in Section 2.3).

If at least one violation is detected, data to be used for the extraction of a new formula are generated by the function `GENAUGMENTEDTRACES` (Line 10). Again, the execution trace $\sigma_{[0,i]}$ is perturbed by adding random Gaussian noise as a counter-overfitting measure, thus producing a set of augmented traces \mathcal{T} of cardinality $n_{aug} + 1$.

Next, function `LEARNFORMULA` (Line 11) tries to extract a set of formulae Φ able to identify a bad behaviour of the system that anticipates, over the traces in \mathcal{T} , the violation de-

tected by the monitor at time instant i . In our case, the extraction is carried out by the genetic algorithm described in Section 4.2, limiting to maximum one formula ϕ , i.e., $|\Phi| \leq 1$. Notice that Φ may be *empty*, an event that occurs if no formula satisfying the quality criteria q can be extracted. The quality criteria relate to the actual ability of a given formula $\phi \in \Phi$ to elicit an anticipatory bad behavior from the traces \mathcal{T} , while maintaining a False Alarm Rate (the fraction of false failure detections, FAR) that does not exceed a specified threshold for the traces in Σ_{\perp} . In setting such a threshold, it should be considered that formulas with a high FAR may cause a degradation of the monitoring pool, where other ill-founded formulas are added as a result of their triggering.

The monitoring executed over all prefixes of σ ends when either σ is correctly recognized as a failure trace by a formula in \mathcal{P} (`break` instruction on Line 13), or σ has run out of points without any failure detection. In the latter case, since trace σ was a failure one, we force the formula extraction process (Lines 15–19). This approach is inspired by the *teacher forcing* technique in deep learning [Williams and Zipser, 1989]. Initially, the framework starts with a possibly empty pool \mathcal{P} of properties. In the extreme case $\mathcal{P} = \emptyset$, it cannot identify any bad behavior of the system and the pool update process is forcibly triggered by the code in Lines 15–19. This is akin to having an oracle guiding the framework. Over time, as \mathcal{P} expands, it is expected to gradually replace the oracle’s function in identifying failures.

During Algorithm 1 operation, the pool of formulas \mathcal{P} is iteratively refined by adding new formulas which ought to predict bad behaviors earlier and with increased reliability and coverage. In addition to this refinement process autonomously operated by the framework, at any time, domain experts can, in principle, make changes to the pool \mathcal{P} , e.g., by manually specifying a new formula encoding a bad behavior.

4.2 Genetic programming algorithm

The function `LEARNFORMULA` is realized by means of a genetic algorithm implementing a genetic programming task, the idea being to evolve formulas starting from an initial population of random solutions [Poli *et al.*, 2008].

In practice, we utilize a multi-objective evolutionary algorithm due to its inherent flexibility. This approach enables us to define a specific grammar that the formulas must adhere to, and it also facilitates the straightforward setup and adjustment of optimization objectives for their extraction. The algorithm is implemented through the library DEAP (Distributed Evolutionary Algorithms in Python) [Fortin *et al.*, 2012].

Figure 2 reports its intuitive operation. The algorithm gets in input the set \mathcal{T} of augmented traces generated by `genAugmentedTraces`, each of length t_l , the set of augmented good training traces Σ_{\perp} , and the set of quality requirements q that the output formula has to satisfy. These include, as we will see, min_{acc} and max_{far} ; they are related to respectively \mathcal{T} and Σ_{\perp} . As for the result, as already mentioned, in our case the algorithm returns a set Φ composed of at most one logic formula such that it satisfies the quality requirements expressed in q and it captures an anticipatory bad behaviour exhibited by traces in \mathcal{T} .

¹Since failures are terminating events, it is reasonable to assume that they can be detected as they occur. This allows to appropriately label the time series and to generate such a training set \mathcal{X} . The same assumption does not hold, e.g., for anomaly detection, a task which is not covered in this paper.

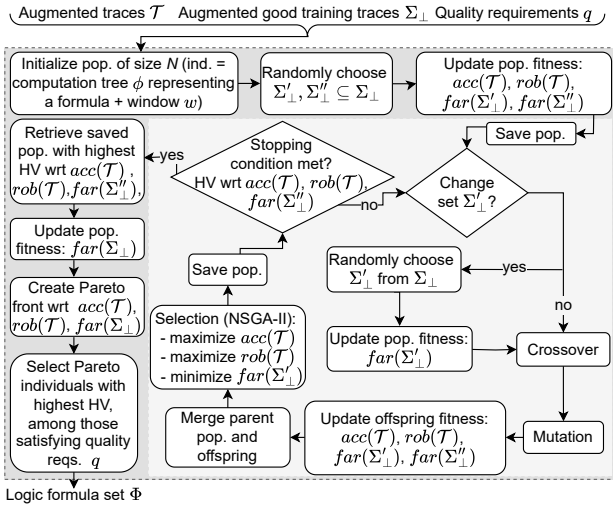


Figure 2: Overview of the genetic algorithm for formula extraction.

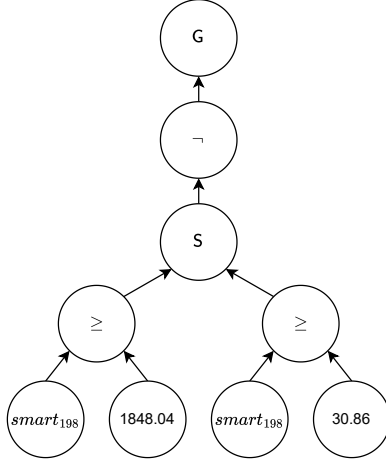


Figure 3: Computation tree of the formula $G(\neg((smart_{198} \geq 1848.04) S (smart_{198} \geq 30.86)))$.

Population and its initialization

As the first step, the population is initialized. Each individual of the population encodes a pair (ϕ, w) , where ϕ is a computation tree representing either a bfSTL or a ppSTL formula (refer to Fig. 3 for an example) and w is an integer in the interval $[1, t_l - 1]$.

The computation tree ϕ is generated using DEAP’s *gen-HalfAndHalf* method, configured to output a tree with a maximum height of 6, as recommended by Koza in his seminal work [Koza, 1994]. Specifically, half the time a tree whose leaves are all at the same depth is returned; in the remaining cases, different leaves may lay at different depths. In addition, care is taken to generate a population of individuals stratified with respect to the height of the trees.

The leaves of a computation tree may represent a variable v that refers to one of the signals that compose each multivariate trace $t \in \mathcal{T}$, or a constant $c \in [0, 1] \subseteq \mathbb{R}$. Note that, although the choice of the domain for the constants may seem

restrictive, in fact, as its first operation, the genetic algorithm normalizes the input traces into the interval $[0, 1]$. Constants c and variables v are compared by means of the \geq operator, in two manners: either a variable is compared with a constant $v \geq c$, or a variable is compared to another variable $v_1 \geq v_2$, where $v_1 \neq v_2$, i.e., they refer to two different signals.

Concerning the internal nodes of the tree, they encode logic operators following the syntax of either bfSTL or ppSTL. In addition, each temporal operator may be (always be, in the case of bfSTL formulas) paired with two interval bounds, i.e., $[a, b]$, with the constants $a, b \in \mathbb{N}$ and $a \leq b$. In generating a computation tree, idempotence is exploited to streamline the formula representation. This is achieved by the definition of a suitable grammar that avoids the combination of redundant operations, such as $\neg\neg\phi$, where ϕ is a bfSTL or a ppSTL formula.

All constants (including the individual’s window w) are implemented through DEAP’s *EphemeralConstant*.

As for w , it is used during the fitness evaluation step to partition each trace $t \in \mathcal{T}$ into a “good behaviour” and a “bad behaviour” sub-trace, as we will now discuss.

Fitness function

Next, the fitness of each individual (ϕ, w) is updated along 4 dimensions. The first 2 pertain to the set of traces \mathcal{T} , while the latter 2 refer, respectively, to the sets Σ_\perp' and Σ_\perp'' . Both are randomly sampled from Σ_\perp , in a way such that $|\Sigma_\perp'| = |\Sigma_\perp''| = \text{fract}_{\text{good}} * |\Sigma_\perp|$ (where $\text{fract}_{\text{good}} \in [0, 1]$ is a global parameter of the framework).

As for the traces in \mathcal{T} , recall that they are all of the same length t_l and their last time instant corresponds to the one immediately preceding the failure detected within the framework training loop, either by means of a formula in the pool \mathcal{P} or through teacher forcing. Thus, we would like the formula ϕ encoded by an individual to be able to elicit an anticipatory bad behaviour of that failure. To such an extent, each trace $t \in \mathcal{T}$ is divided into a good behaviour sub-trace $\text{good}(t) = t_{[0, w-1]}$ and a bad behaviour sub-trace $\text{bad}(t)$ as follows:

$$\text{bad}(t) = \begin{cases} t_{[w, t_l-1]} & \text{if formula } \phi \text{ is bfSTL} \\ t_{[0, w]} & \text{if formula } \phi \text{ is ppSTL} \end{cases}$$

The idea is that the last time points of t are those closer to the detected failure, and thus should contain some prelude of it.

By means of `rtamt`, the monitor mon_ϕ is applied over all good and bad sub-traces, and the first fitness element, to be maximized, is obtained as follows:

$$\text{acc}(\mathcal{T}) = \frac{|\{t \mid t \in \mathcal{T} \wedge \text{mon}_{G\phi}(\text{good}(t)) \in \{\top, ?\}\}|}{|\mathcal{T}|} + \frac{|\{t \mid t \in \mathcal{T} \wedge \text{mon}_\phi(\text{bad}(t)) = \perp\}|}{|\mathcal{T}|}.$$

Intuitively, $\text{acc}(\mathcal{T})$ represents the fraction of sub-traces correctly identified as good and failure ones; note how formula $G\phi$ is considered on $\text{good}(t)$, as ϕ must be checked with respect to every time point in such traces.

The second fitness element, to be maximized, refers to the quantitative semantics of STL [Donzé and Maler, 2010]:

$$\text{rob}(\mathcal{T}) = \min(\text{rob}_{\text{good}}, -\text{rob}_{\text{bad}}),$$

where rob_{good} (resp., rob_{bad}) is the average *robustness* of mon_ϕ calculated over all good (resp., bad) sub-traces. Due to the traces normalization step performed at the beginning of the genetic algorithm, in our case, both are real numbers in the interval $[-1, 1]$. Intuitively, the robustness of a formula ϕ with respect to a trace t tells us how far (or close) the formula is to be satisfied on a trace. For example, given $t_1 = [1, 0.9, 0.7]$, $t_2 = [0, 0.3, 0.4]$, and $\phi = G(x \geq 0.5)$, the former trace has a robustness of 0.2, while the latter of -0.1 . Note how the sign is switched on rob_{bad} to make it positive. Maximizing such a formula is akin to the concept of maximum margin separating hyperplane in the context of support vector machines [Hastie *et al.*, 2009].

As for set Σ'_\perp , again, $rtamt$ is applied over all its traces, and the related third fitness element, to be minimized, intuitively represents the fraction of good traces of Σ'_\perp in which a failure is incorrectly signalled. It is obtained as:

$$far(\Sigma'_\perp) = \frac{|\{\sigma \mid \sigma \in \Sigma'_\perp \wedge mon_{G\phi}(\sigma) = \perp\}|}{|\Sigma'_\perp|}.$$

Similarly, for the fitness element over Σ''_\perp , we calculate the fraction of good traces of Σ''_\perp in which a failure is incorrectly signalled:

$$far(\Sigma''_\perp) = \frac{|\{\sigma \mid \sigma \in \Sigma''_\perp \wedge mon_{G\phi}(\sigma) = \perp\}|}{|\Sigma''_\perp|}.$$

Objectives $acc(\mathcal{T})$, $rob(\mathcal{T})$ and $far(\Sigma'_\perp)$ guide the evolutionary process, while $far(\Sigma''_\perp)$ is used to implement an early stopping strategy (which requires a fixed, i.e., not changing during the generations, reference set).

The initialization step of the algorithm finishes by saving the population along with its four fitness elements and the hypervolume calculated with respect to $acc(\mathcal{T})$, $rob(\mathcal{T})$ and $far(\Sigma'_\perp)$.

Let us now focus on the evolutionary part of the algorithm. Every $r_{interval}$ generations, the set Σ'_\perp is resampled from Σ_\perp as a counter-overfitting measure (see, for instance, [Gonçalves and Silva, 2013]). If a resample is performed, the fitness element $far(\Sigma'_\perp)$ is updated for each individual. Then, crossover and mutation operators are applied over the population, generating an offspring. Such operators are defined as follows.

Crossover

Let $i_1 = (\phi_1, w_1)$, $i_2 = (\phi_2, w_2)$ be two individuals from the population. From them, two other individuals can be created according to several crossover operations, randomly selected with uniform probability. The simplest one involves exchanging their windows, thus obtaining two new individuals $i'_1 = (\phi_1, w_2)$, $i'_2 = (\phi_2, w_1)$.

Also, the two computation trees corresponding to ϕ_1 and ϕ_2 can be hybridized by means of DEAP's *cxOnePointLeaf-Biased* operator. In short, it randomly selects a crossover point in each tree and exchanges the subtrees rooted in them. Again following suggestions by Koza [Koza, 1994], our operator is biased to choose the crossover point on internal nodes 90% of the times, while a leaf is chosen 10% of the times. Finally, we set a limit of 17 [Koza, 1994] on the generated

individuals' height (DEAP's *staticLimit*). When an over-the-height-limit individual is generated, it is simply replaced by one of its parents, randomly selected.

Mutation

Given an individual $i = (\phi, w)$, also the mutation can be performed according to several strategies, chosen according to a uniform random probability. As for the window w , it can be re-generated by randomly choosing a number in the interval $[1, t_l - 1]$; or, it can be adjusted by a small random amount of points to the left or right, to allow for a finer tuning.

Concerning the computation tree corresponding to ϕ , three operations can be performed, taken from the DEAP primitives: *mutNodeReplacement*, which replaces a randomly chosen node by a randomly chosen operation with the same number of arguments; *mutShrink*, which shrinks the tree by choosing randomly a branch and replacing it with one of the branch's arguments (also randomly chosen); and, *mutEphemeral*, which randomly changes the values of all of the constants in the tree.

Once the offspring has been produced, their four previously defined fitness elements are determined. Next, the parent population and the offspring are merged, and a selection is performed relying on the classic strategy implemented in NSGA-II [Deb *et al.*, 2002], based on the concepts of *ranking* and *crowding distance*. Specifically, for each individual, the fitness elements $acc(\mathcal{T})$ (maximized), $rob(\mathcal{T})$ (maximized) and $far(\Sigma'_\perp)$ (minimized) are considered. The newly obtained population is then saved, along with its fitness elements and the hypervolume calculated with respect to $acc(\mathcal{T})$, $rob(\mathcal{T})$, and $far(\Sigma''_\perp)$.

The evolutionary part may end according to two conditions: either the maximum number of generations max_{gen} (global parameter of the framework) has been reached, or the early stopping condition is triggered. The latter is defined as follows: the population hypervolume, calculated with respect to $acc(\mathcal{T})$, $rob(\mathcal{T})$ and $far(\Sigma'_\perp)$ (i.e., the fixed set of good traces), is tracked along the generations. Then, if no improvement is observed for *patience* generations, the execution is halted.

Recall that, during the execution of the algorithm, each generation's population has been saved. After the end of the iterative part of the evolutionary process, the population with the highest hypervolume with respect to $acc(\mathcal{T})$, $rob(\mathcal{T})$ and $far(\Sigma'_\perp)$ is recovered, and its individuals' false alarm rate with respect to all Σ_\perp traces, $far(\Sigma_\perp)$, is established.² Then, the Pareto front of optimal solutions with respect to $acc(\mathcal{T})$, $rob(\mathcal{T})$ and $far(\Sigma_\perp)$ is extracted and filtered to keep only individuals satisfying the quality criteria q : $acc(\mathcal{T}) > min_{acc}$ and $far(\Sigma_\perp) \leq max_{far}$, where $min_{acc}, max_{far} \in [0, 1]$. Also, they are required to correctly classify the good and bad sub-traces originated from the original, unperturbed trace, always included in \mathcal{T} . Finally, among the remaining individuals, the formula of the best performing one is returned. This latter selection is based on choosing the individual with the highest hypervolume, calculated in relation to $acc(\mathcal{T})$, $rob(\mathcal{T})$, and $far(\Sigma_\perp)$.

²We use Σ''_\perp instead of Σ_\perp for computational efficiency.

4.3 Comparison to related work

Let us briefly contrast our framework with existing approaches to the integration of monitoring and learning.

Compared to contributions in the literature that extract temporal relations in time series data [Aggarwal *et al.*, 2018; Lu *et al.*, 2020; Petmezas *et al.*, 2021; Gao *et al.*, 2022], which are based on black-box techniques such as deep learning, our framework is highly interpretable (we refer the reader to Appendix E for some examples of extracted formulas). Interpretability is a paramount requirement in critical scenarios, such as healthcare and avionics. This is relevant especially in the light of recent results that underscore the issues behind feature attribution methods [Bilodeau *et al.*, 2024].

The closest proposal is the one described in [Brunello *et al.*, 2023], that uses genetic programming to learn new formulas of bfSTL. Here, we improved it under several respects.

From the theoretical point of view, we introduced the unbounded past variant of STL and defined the formalisms G(ppSTL) and F(ppSTL) which are strictly more expressive than bfSTL, that is, they can express more properties while still remaining in the realm of monitorability, removing the need of checking whether a newly learned formula is monitorable. In such respect, recall that, while bfSTL can only constrain bounded intervals of a trace, G(ppSTL) and F(ppSTL) do not inherit such a limitation, being able to constrain arbitrarily long intervals. Extending the learning algorithm of [Brunello *et al.*, 2023] to the new formalism was not trivial, and it advantageously led to the dismissal of the concept of *horizon*. The generation of monitorable-only formulas in the genetic algorithm is ensured by the design of a type-based grammar that allows only well-formed formulas adhering to the chosen formalism to be generated, while removing redundant operators (exploiting idempotence). In addition to these changes necessary to deal with the past, we introduced several new features:

- quality requirements for the new formulas are now clearly stated by means of q . Above all, it can be ensured that a formula’s FAR (*i.e.*, the False Alarm Rate) does not exceed a specified threshold. This was not possible previously, resulting, over time, in a large number of formulas being removed from the pool after they had been learned;
- the imposed FAR threshold allowed us to redesign the offline learning phase of the framework increasing its efficiency. Now, it is sufficient to iterate only over failures, rather than the entire dataset;
- the refactored genetic algorithm makes use of a rotation-based strategy to counter overfitting on the good traces, includes an enriched set of crossover and mutation operations, and exploits an enhanced early stopping strategy, that returns the best population among the observed ones.

5 Experimental evaluation

Here, we present the results of the application of the framework on three well-known datasets from the literature, comparing its performance with other recent contributions.

The *Backblaze Hard Drive* dataset contains information on the “health” status of hard drives, tracked by means of Self Monitoring Analysis and Reporting Technology (SMART). Each trace is described by: date of the report, serial number of the drive, a label indicating a drive failure, and 21 numerical SMART parameters. For comparison with the literature, we focus on [Brunello *et al.*, 2023] Split *S1*. The *Tennessee Eastman Process* (TEP) dataset is composed of simulated data from a fictitious chemical plant. Each trace has the features: trace ID, normal/faulty label, and 52 variables tracking data about the operating values of plant components. The *NASA Commercial Modular Aero-Propulsion System Simulation* (C-MAPSS) dataset includes run-to-failure simulated data of turbofan jet engines. Specifically, we focus on the dataset FD001 and on the detection of the *Unhealthy state* class [Kim and Sohn, 2020]. Each simulation is represented by a multivariate time series, sampled at one value per second, obtained from 21 engine sensors. Details about the three datasets are included in Appendix C.

We instantiate two versions of the framework, using two different fragments of STL. The first, GP-bfSTL, relies on bfSTL formulas, as specified in the work by [Brunello *et al.*, 2023]. The second, GP-G(ppSTL), uses the more expressive G(ppSTL) fragment. Thereafter, we proceeded by first confirming the literature results by means of GP-bfSTL. Then, given such a baseline, we assess the performance of GP-G(ppSTL) to determine the contribution brought by the newly isolated logical fragment.

For each dataset, the GP-bfSTL framework is tuned considering a separate validation set obtained from training data (details are provided in Appendix D). Then, the best hyper parameters are used to run both the GP-bfSTL and GP-G(ppSTL) framework over all the training sets. Finally, performance is established by applying the monitor to classify each test set trace independently by means of the pool of properties obtained during training. To account for the inherent stochasticity of our approach, each experiment is repeated various times (of course, the same seeds are used with GP-bfSTL and GP-G(ppSTL)). Appendix F reports all the details about the code and how to use it to reproduce and use our framework for new research.

5.1 Results

Table 1 reports the main outcomes of the experiments; a discussion about the metrics employed as well as additional findings and analyses are in Appendix E. Note that the competitors’ results are based on a single execution (still reported under the Avg columns), except for [Brunello *et al.*, 2023]. Thus, to allow for a fair comparison, for the F1 and MCC³ metrics, which summarize the overall performance of the framework, we report not only the average but also the minimum and maximum values observed across the repetitions.

Overall, based on the F1 metric, which is provided by all the competitors, we can observe that GP-bfSTL achieved results on par with or superior to the considered baselines, confirming the soundness of our implementation. The GP-

³Matthews Correlation Coefficient, which is more informative than F1 in several scenarios [Chicco and Jurman, 2020].

Table 1: Experimental results

Dataset	Approach	Prec.	Rec.	FAR	Min	FI	Avg	Min	MCC	Avg
Backblaze S1	[Huang, 2017]	.51	.54	.00	-	-	.52	-	-	-
	[Lu <i>et al.</i> , 2020]	.87	.41	.00	-	-	.55	-	-	-
	[Brunello <i>et al.</i> , 2023]	.54	.60	.00	-	-	.56	-	-	-
	GP-bfSTL (our)	.76	.49	.00	.58	.61	.59	.58	.63	.61
	GP-G(ppSTL) (our)	.85	.47	.00	.55	.62	.60	.59	.64	.62
TEP	[Hajihosseini <i>et al.</i> , 2018]	1.0	1.0	-	-	-	1.0	-	-	-
	[Onel <i>et al.</i> , 2019]	1.0	1.0	.00	-	-	1.0	-	-	-
	[Brunello <i>et al.</i> , 2023]	1.0	1.0	.00	-	-	1.0	-	-	-
	GP-bfSTL (our)	1.0	1.0	.00	1.0	1.0	1.0	1.0	1.0	1.0
	GP-G(ppSTL) (our)	1.0	1.0	.00	1.0	1.0	1.0	.99	1.0	1.0
C-MAPSS	[Kim and Sohn, 2020]	.71	1.0	-	-	-	.83	-	-	-
	[Brunello <i>et al.</i> , 2023]	.96*	.77*	.01*	-	-	.86*	-	-	-
	GP-bfSTL (our)	.95	.65	.02	.70	.85	.77	.61	.80	.69
	GP-G(ppSTL) (our)	.98	.67	.01	.71	.87	.79	.63	.82	.73

Note: Results of [Lu *et al.*, 2020] listed as in [Brunello *et al.*, 2023]; the others are reported as in the original references. *: the original results in the referenced paper were based on an incorrectly generated test set of 143 traces, rather than the official set of 100 traces. The latter has been correctly considered in the present work.

G(ppSTL) version, when run with the same hyper parameters as the GP-bfSTL one, scored in par or even better (cf. MCC), confirming the theoretical claims of Section 3. It is plausible that a dedicated tuning phase, which could not be performed due to time/resource limits, could enhance results further.

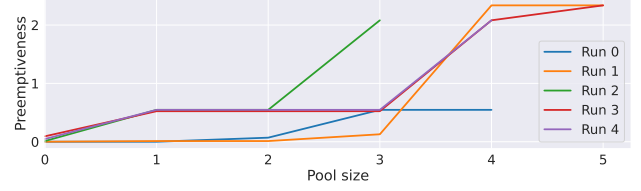
Note that a significant gap can sometimes be observed between the maximum/minimum performances of our framework. This stems from our chosen experimental workflow: for statistical soundness, also the framework’s iterative learning process relies on a random shuffle of traces (Algorithm 1, Line 4), which in turn affects the formula extraction process. We speculate that learning on the training set through multiple passes until the pool converges (i.e., when no more properties are added) would yield more stable results.

Focusing on the GP-G(ppSTL) version of the framework, being based on the newly isolated fragment and the best performing one, Figure 4 illustrates how the average preemptiveness in failure detection on test set instances evolves as more and more formulas are learned and incorporated into the monitoring pool. In the dataset TEP, a plateau is reached around a value of 200, meaning that the formulas in the pool are capable of (perfectly, given Table 1) identifying a failure roughly 10 hours before it happens (recall the 3-minute sampling time used for the time series in this dataset). An increase in the anticipatory behaviour of the framework is observed also in the datasets C-MAPSS and Backblaze S1.

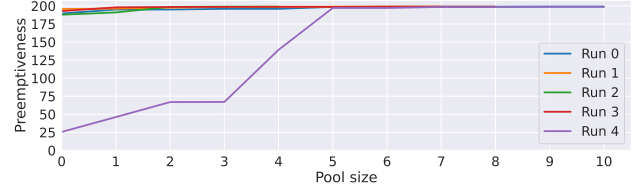
Finally, we report an example of formulas extracted during the framework training phase for the Backblaze S1 dataset.

Formula $f_1 = G(\text{smart}_{187} < 24.40 \vee \text{smart}_1 \geq 169483444)$, included in \mathcal{P} , evaluates to false, thus issuing a failure detection. It means that the number of *reported uncorrectable errors* (sensor smart_{187}) exceeded the 24.40 threshold and, at the same time, the *read error rate* (sensor smart_1) was lower than 169483444. As a result, another formula, intended to anticipate the situation modeled by f_1 , is extracted by the genetic algorithm and added to \mathcal{P} . It is $f_2 = G(\text{smart}_{198} < 16.56)$, which is violated when the *uncorrectable sector count* exceeds the 16.56 threshold.

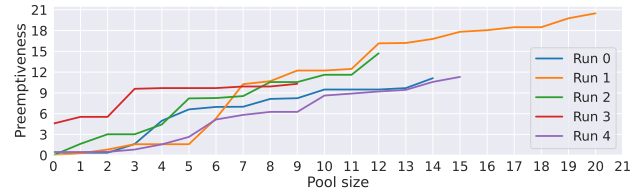
Thus, when viewed from the perspective of the domain, an event marked by the violation of f_2 can be regarded as the direct cause of the scenario projected by the violation of f_1 . In other words, hardware errors (indicated by f_2 ’s violation) are the root cause of subsequent read/write errors (as modeled by the violation of f_1).



(a) Backblaze S1 dataset (preemptiveness in days).



(b) TEP dataset (preemptiveness in 3-minute steps).



(c) C-MAPSS dataset (preemptiveness in seconds).

Figure 4: Average preemptiveness evolution on the test set instances, as more properties are learned and added to the pool, for 5 different runs of the framework, GP-G(ppSTL) case.

6 Conclusions and discussion

In this paper, we propose to exploit machine learning to enhance monitoring of STL formulas. Our methodology is based on the identification of two fragments of STL, namely G(ppSTL) and F(ppSTL), for which we give formal guarantees on monitorability and expressiveness, and on the use of genetic programming to automatically learn new relevant formulas of these fragments. The experiments reveal that formulas generated in this way are effective in anticipating the discovery of failures.

We conclude by discussing future research directions, in particular about formula extraction and management. The use of a genetic algorithm facilitated rapid prototyping and offered flexibility. Still, more efficient methodologies should be explored. For example, the use of generative and/or reinforcement learning techniques [Holt *et al.*, 2023], or a hybrid approach combining deep learning with evolutionary computation, could offer substantial advancements [Chen *et al.*, 2018; Chen *et al.*, 2020; Mundhenk *et al.*, 2021; Qian *et al.*, 2021]. Graph neural networks are promising as well, as they are proven effective for performing symbolic regression tasks [Cranmer *et al.*, 2020], and may more effectively exploit formulas, e.g., considering their automaton or directed acyclic graph encodings. Finally, the framework itself should be expanded to support online continual learning and update of formulas, rather than relying solely on a fixed training dataset. This is crucial for real-world applications where monitored systems may change their behaviour over time.

Acknowledgments

All the authors acknowledge the support from the 2024 Italian INdAM-GNCS project “Certificazione, monitoraggio, ed interpretabilità in sistemi di intelligenza artificiale”, ref. no. CUP E53C23001670001. Luca Geatti, Angelo Montanari, and Nicola Saccomanno also acknowledge the support from the Interconnected Nord-Est Innovation Ecosystem (iNEST), which received funding from the European Union Next-GenerationEU (PIANO NAZIONALE DI RIPRESA E RESILIENZA (PNRR) – MISSIONE 4 COMPONENTE 2, INVESTIMENTO 1.5 – D.D. 1058 23/06/2022, ECS00000043). In addition, Angelo Montanari acknowledges the support from the MUR PNRR project FAIR - Future AI Research (PE00000013) also funded by the European Union Next-GenerationEU. This manuscript reflects only the authors’ views and opinions, neither the European Union nor the European Commission can be considered responsible for them.

References

- [Aggarwal *et al.*, 2018] Karan Aggarwal, Onur Atan, Ahmed K Farahat, Chi Zhang, Kosta Ristovski, and Chetan Gupta. Two birds with one network: Unifying failure event prediction and time-to-failure modeling. In *Proceedings of IEEE International Conference on Big Data (Big Data)*, pages 1308–1317. IEEE, 2018.
- [Bauer *et al.*, 2011] Andreas Bauer, Martin Leucker, and Christian Schallhart. Runtime verification for LTL and TLTL. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 20(4):1–64, 2011.
- [Bilodeau *et al.*, 2024] Blair Bilodeau, Natasha Jaques, Pang Wei Koh, and Been Kim. Impossibility theorems for feature attribution. *Proceedings of the National Academy of Sciences*, 121(2):e2304406120, 2024.
- [Brunello *et al.*, 2023] Andrea Brunello, Dario Della Monica, Angelo Montanari, Nicola Saccomanno, and Andrea Urgolo. Monitors that learn from failures: Pairing STL and genetic programming. *IEEE Access*, 11:57349–57364, 2023.
- [Chang *et al.*, 1992] Edward Y. Chang, Zohar Manna, and Amir Pnueli. Characterization of temporal property classes. In *Proceedings of the 19th International Colloquium on Automata, Languages and Programming (ICALP 1992)*, pages 474–486, 1992.
- [Chen *et al.*, 2018] Yukang Chen, Gaofeng Meng, Qian Zhang, Shiming Xiang, Chang Huang, Lisen Mu, and Xinggang Wang. Reinforced evolutionary neural architecture search. *arXiv preprint arXiv:1808.00193*, 2018.
- [Chen *et al.*, 2020] Diqi Chen, Yizhou Wang, and Wen Gao. Combining a gradient-based method and an evolution strategy for multi-objective reinforcement learning. *Applied Intelligence*, 50:3301–3317, 2020.
- [Chicco and Jurman, 2020] Davide Chicco and Giuseppe Jurman. The advantages of the matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics*, 21(1):1–13, 2020.
- [Cranmer *et al.*, 2020] M. Cranmer, Alvaro Sanchez-Gonzalez, P. Battaglia, Rui Xu, K. Cranmer, D. Spergel, and S. Ho. Discovering symbolic models from deep learning with inductive biases. *ArXiv*, abs/2006.11287, 2020.
- [Deb *et al.*, 2002] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computation*, 6(2):182–197, 2002.
- [Donzé and Maler, 2010] Alexandre Donzé and Oded Maler. Robust satisfaction of temporal logic over real-valued signals. In Krishnendu Chatterjee and Thomas A. Henzinger, editors, *Proceedings of the 8th International Conference on Formal Modeling and Analysis of Timed Systems (FORMATS 2010)*, pages 92–106. Springer, 2010.
- [Fortin *et al.*, 2012] F. A. Fortin, F. M. De Rainville, M. A. Gardner, M. Parizeau, and C. Gagné. DEAP: Evolutionary algorithms made easy. *Journal of Machine Learning Research*, 13:2171–2175, 2012.
- [Gao *et al.*, 2022] Jiechao Gao, Haoyu Wang, and Haiying Shen. Task failure prediction in cloud data centers using deep learning. *IEEE Transactions on Services Computing*, 15:1411–1422, 2022.
- [Ghallab *et al.*, 2004] Malik Ghallab, Dana S. Nau, and Paolo Traverso. *Automated planning - theory and practice*. Elsevier, 2004.
- [Gonçalves and Silva, 2013] Ivo Gonçalves and Sara Silva. Balancing learning and overfitting in genetic programming with interleaved sampling of training data. In *Proceedings of the European Conference on Genetic Programming (EuroGP)*, pages 73–84. Springer, 2013.
- [Hajihosseini *et al.*, 2018] Payman Hajihosseini, Mohammad Mousavi Anzehaee, and Behzad Behnam. Fault detection and isolation in the challenging Tennessee Eastman Process by using image processing techniques. *ISA Transactions*, 79:137–146, 2018.
- [Hastie *et al.*, 2009] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- [Havelund and Peled, 2023] Klaus Havelund and Doron Peled. Monitorability for runtime verification. In Panagiotis Katsaros and Laura Nenzi, editors, *Proceedings of the 23rd International Conference on Runtime Verification (RV 2023)*, volume 14245 of *Lecture Notes in Computer Science*, pages 447–460. Springer, 2023.
- [Holt *et al.*, 2023] Samuel Holt, Zhaozhi Qian, and Mihaela van der Schaar. Deep generative symbolic regression. *arXiv preprint arXiv:2401.00282*, 2023.
- [Huang, 2017] X. Huang. Hard drive failure prediction for large scale storage system. 2017. M.Sc Thesis.
- [Kim and Sohn, 2020] Tae San Kim and So Young Sohn. Multitask learning for health condition identification and

- remaining useful life prediction: Deep convolutional neural network approach. *Journal of Intelligent Manufacturing*, 32(8):2169–2179, 2020.
- [Koza, 1994] John R Koza. Genetic programming as a means for programming computers by natural selection. *Statistics and Computing*, 4(2):87–112, 1994.
- [Kupferman and Vardi, 2001] Orna Kupferman and Moshe Y Vardi. Model checking of safety properties. *Formal Methods in System Design*, 19(3):291–314, 2001.
- [Leucker and Schallhart, 2009] Martin Leucker and Christian Schallhart. A brief account of runtime verification. *The Journal of Logic and Algebraic Programming*, 78(5):293–303, 2009.
- [Lu et al., 2020] Sidi Lu, Bing Luo, Tirthak Patel, Yongtao Yao, Devesh Tiwari, and Weisong Shi. Making disk failure predictions SMARTer! In *Proceedings of the 21st USENIX Conference on File and Storage Technologies (FAST’23)*, pages 151–167. USENIX Association, 2020.
- [Maler and Nickovic, 2004] Oded Maler and Dejan Nickovic. Monitoring temporal properties of continuous signals. In Yassine Lakhnech and Sergio Yovine, editors, *Proceedings of the Joint International Conferences on Formal Modelling and Analysis of Timed Systems (FORMATS 2004) and Formal Techniques in Real-Time and Fault-Tolerant Systems (FTRTFT 2004)*, pages 152–166. Springer, 2004.
- [Maler et al., 2007] Oded Maler, Dejan Nickovic, and Amir Pnueli. On synthesizing controllers from bounded-response properties. In *Proceedings of the 19th International Conference on Computer Aided Verification (CAV 2007)*, pages 95–107. Springer, 2007.
- [Mundhenk et al., 2021] Terrell Mundhenk, Mikel Landajuela, Ruben Glatt, Claudio P Santiago, Brenden K Petersen, et al. Symbolic regression via deep reinforcement learning enhanced genetic programming seeding. *Advances in Neural Information Processing Systems*, 34:24912–24923, 2021.
- [Ničković and Yamaguchi, 2020] D. Ničković and T. Yamaguchi. RTAMT: Online robustness monitors from STL. In *Proceedings of the 18th International Symposium on Automated Technology for Verification and Analysis (ATVA 2020)*, volume 12302, pages 564–571. Springer, 2020.
- [Onel et al., 2019] Melis Onel, Chris A. Kieslich, and Efstratios N. Pistikopoulos. A nonlinear support vector machine-based feature selection approach for fault detection and diagnosis: Application to the Tennessee Eastman Process. *American Institute of Chemical Engineers Journal*, 65(3):992–1005, 2019.
- [Petmezas et al., 2021] Georgios Petmezas, Kostas Haris, Leandros Stefanopoulos, Vassilis Kilintzis, Andreas Tzavelis, John A Rogers, Aggelos K Katsaggelos, and Nicos Maglaveras. Automated atrial fibrillation detection using a hybrid CNN-LSTM network on imbalanced ECG datasets. *Biomedical Signal Processing and Control*, 63:102194, 2021.
- [Poli et al., 2008] R. Poli, W. Langdon, and N. Mcphee. *A Field Guide to Genetic Programming*. <http://www.gp-field-guide.org.uk>, 2008.
- [Qian et al., 2021] Kai Qian, Jie Jiang, Yulong Ding, and Shuang-Hua Yang. DLGEA: A deep learning guided evolutionary algorithm for water contamination source identification. *Neural Computing and Applications*, 33:11889–11903, 2021.
- [Williams and Zipser, 1989] Ronald J. Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*, 1(2):270–280, 1989.

A Proofs of theorems

Proposition 2 (ppSTL is more expressive than bfSTL). *There exists a language $\mathcal{L} \subseteq (\mathbb{T} \times \mathbb{D}^n)^*$ (for some $n \in \mathbb{N}$) such that:*

- *there exists a ppSTL formula ϕ such that $\mathcal{L}(\phi) = \mathcal{L}$;*
- *there exists no bfSTL formula ψ such that $\mathcal{L}(\psi) = \mathcal{L}$.*

Proof. Let $n = 1$ and let \mathcal{L} be the set of real-time 1-valued traces $\sigma = \langle (t_0, v_0), \dots, (t_n, v_n) \rangle$ such that $v_0 \geq 0$ iff $v_n \geq 0$, that is the set of traces on which the value of the variable x_0 at the *first* position is greater than 0 iff the value of x_0 at the *last* position is greater than 0. The ppSTL formula $\phi := (x_0 \geq 0) \longleftrightarrow \text{O}(\neg \text{Y}(\perp) \wedge (x_0 \geq 0))$ is such that $\mathcal{L}(\phi) = \mathcal{L}$.

To prove that there exists no bfSTL formula ψ such that $\mathcal{L}(\psi) = \mathcal{L}$, we proceed by contradiction. Let ψ be a formula of bfSTL such that $\mathcal{L}(\psi) = \mathcal{L}$. By definition of the bounded future fragment of STL, there exists a number $d \in \mathbb{N}$ (which is computable starting from the formula ϕ and is called the *temporal depth* of ϕ , cf. [Maler et al., 2007]) such that, for all $\sigma' \in (\mathbb{T} \times \mathbb{D}^n)^*$, it holds that $\sigma_{[0,d]} \cdot \sigma' \models \psi$ iff $\sigma_{[0,d]} \models \psi$. In other words, the satisfaction of ψ by σ depends only on the prefix of σ from 0 to d , where d is a number computable from ψ . Consider now any trace σ of length more than $d + 1$ such that:

- the value of the variable x_0 at the *first* position is greater than 0 iff the value of x_0 at the *last* position is greater than 0;
- the value of the variable x_0 at the *first* position is greater than 0 iff the value of x_0 at position d is *not* greater than 0;

Clearly, $\sigma \in \mathcal{L}$. However, $\sigma_{[0,d]} \notin \mathcal{L}$. Since ψ is a bfSTL formula, it follows that $\sigma_{[0,d]} \cdot \sigma' \notin \mathcal{L}(\psi)$, for all $\sigma' \in (\mathbb{T} \times \mathbb{D}^n)^*$. In particular, $\sigma \notin \mathcal{L}(\psi)$. This is a contradiction with the fact that $\mathcal{L}(\psi) = \mathcal{L}$. \square

Lemma 1. *For all $\phi \in \text{G}(\text{ppSTL})$ (resp., $\phi \in \text{F}(\text{ppSTL})$), it holds that $\mathcal{L}(\phi)$ is a safety (resp., cosafety) real-time multi-valued language.*

Proof. We begin with the case of $\text{F}(\text{ppSTL})$. Let ϕ be a formula of $\text{F}(\text{ppSTL})$ over n variables. By Definition 7, it holds that $\phi = \text{F}(\psi)$ for some $\psi \in \text{ppSTL}$. Let $\sigma \in (\mathbb{T} \times \mathbb{D}^n)^\omega$ be a trace such that $\sigma \in \mathcal{L}(\phi)$. By the semantics of the F operator, it holds that there exists an $i \geq 0$ such that $\sigma_{[0,i]} \models \psi$. Since ψ is a pure past formula of STL, it holds that $\sigma_{[0,i]} \cdot \sigma', i \models \psi$ for all $\sigma' \in (\mathbb{T} \times \mathbb{D}^n)^\omega$. Therefore, $\sigma_{[0,i]} \cdot \sigma' \in \mathcal{L}(\phi)$ for all $\sigma' \in (\mathbb{T} \times \mathbb{D}^n)^\omega$. By Definition 3, it follows that $\mathcal{L}(\phi)$ is cosafety. We now prove the case of $\text{G}(\text{ppSTL})$. Let $\phi \in \text{G}(\text{ppSTL})$. By definition of this fragment, we have that $\phi = \text{G}(\psi)$ for some $\psi \in \text{ppSTL}$. Let $\phi' := \text{F}(\neg\psi)$. It holds that $\mathcal{L}(\phi) = \overline{\mathcal{L}(\phi')}$. Moreover, since $\phi' \in \text{F}(\text{ppSTL})$, we have that $\mathcal{L}(\phi')$ is cosafety. It follows from Definition 4 that $\mathcal{L}(\phi)$ is safety. \square

A.1 Completeness for the qualitative-time case

In this part, we focus on the *qualitative-time* case and we prove that, under this condition, the fragments $\text{G}(\text{ppSTL})$ and $\text{F}(\text{ppSTL})$ are expressively *complete*, that is, every safety

(resp., cosafety) multi-valued language that can be defined by an STL formula is definable also in $\text{G}(\text{ppSTL})$ (resp., in $\text{F}(\text{ppSTL})$).

To prove the result, we first define the notion of Boolean abstraction of an STL formula. From now on, for each formula ϕ in STL, we let \mathcal{AF}_ϕ be the set of atomic subformulas in ϕ (that is, those of type $f_i(\bar{x}) \otimes d$).

Definition 8 (Boolean Abstraction of STL formulas). *Let ϕ be a formula of STL. We define the Boolean abstraction of ϕ , denoted as $\mathbb{B}(\phi)$, as the MTL formula $\phi' \wedge C$ over \mathcal{AP} where:*

- $\mathcal{AP} := \{p_\psi \mid \psi \in \mathcal{AF}_\phi\}$
- ϕ' is the formula obtained from ϕ by replacing each atomic subformula in \mathcal{AF}_ϕ with the atomic proposition $p_\psi \in \mathcal{AP}$;
- C is the conjunction of the formulas $\text{G}(\neg p_\psi \vee \neg p_{\psi'})$, for all subformulas $\psi, \psi' \in \mathcal{AF}_\psi$ such that $\psi \wedge \psi'$ is unsatisfiable.

Similarly, for any MTL formula γ over $\mathcal{AP} = \{p_\psi \mid \psi \in \mathcal{AF}_\phi\}$, we define $\mathbb{B}^{-1}(\gamma)$ as the formula obtained from γ by replacing each $p_\psi \in \mathcal{AP}$ with $\psi \in \mathcal{AF}_\phi$.

As an example, consider the formula $\text{G}(\text{O}\psi \wedge \text{O}\psi')$, where $\psi := x > 0$ and $\psi' := y < 0 \wedge x < -5$. The Boolean abstraction in this case is the following MTL formula:

$$\text{G}(\text{O}p_\psi \wedge \text{O}p_{\psi'}) \wedge \text{G}(\neg p_\psi \vee \neg p_{\psi'})$$

The following lemma is crucial for the completeness theorem: it shows that the application of $\mathbb{B}(\cdot)$ preserves the safety condition.

Lemma 2. *For each $\phi \in \text{STL}$, if $\mathcal{L}(\phi)$ is safety (resp., cosafety) then $\mathcal{L}(\mathbb{B}(\phi))$ is safety (resp., cosafety).*

Proof. We give the proof for the safety case (the cosafety case is identical). Let ϕ be a formula in STL, and let $\mathcal{AP} := \{p_\psi \mid \psi \in \mathcal{AF}_\phi\}$. We define the bijection $\nu : (\mathbb{N} \times \mathbb{D}^n)^\omega \rightarrow (\mathbb{N} \times 2^{\mathcal{AP}})^\omega$ as follows: $\nu(\langle \sigma_0, \sigma_1, \dots \rangle) := \langle \sigma'_0, \sigma'_1, \dots \rangle$ such that $\sigma'_i := \{p_\psi \mid \sigma_i \models \psi, \forall \psi \in \mathcal{AF}_\phi\}$, for each $i \geq 0$. By a straightforward induction on the structure of ϕ , we can prove that $\sigma \in \mathcal{L}(\phi)$ iff $\nu(\sigma) \in \mathcal{L}(\mathbb{B}(\phi))$.

Now, *ad absurdum* we suppose that $\mathcal{L}(\mathbb{B}(\phi))$ is *not* a safety language. By Definition 4, this means that there exists a trace $\sigma \notin \mathcal{L}(\mathbb{B}(\phi))$ such that, for all $i \geq 0$ there exists a $\sigma' \in (\mathbb{N} \times 2^{\mathcal{AP}})^\omega$ such that $\sigma_{[0,i]} \cdot \sigma' \in \mathcal{L}(\mathbb{B}(\phi))$. Consider now the trace $\nu^{-1}(\sigma)$. It holds that $\nu^{-1}(\sigma) \notin \mathcal{L}(\phi)$ and, for all $i \geq 0$, there exists a $\sigma' \in (\mathbb{N} \times 2^{\mathcal{AP}})^\omega$ such that $\nu^{-1}(\sigma_{[0,i]} \cdot \sigma') \in \mathcal{L}(\phi)$. This implies that $\mathcal{L}(\phi)$ is *not* a safety language, but this is a contradiction with the hypothesis. Therefore, $\mathcal{L}(\mathbb{B}(\phi))$ has to be a safety language. \square

We are now ready to prove that, in the case of qualitative time, all safety multi-valued languages definable in STL are also definable in $\text{G}(\text{ppSTL})$, and that the same holds for the cosafety case.

Theorem 2 (Expressive Completeness over qualitative-time). *For all multi-valued languages $\mathcal{L} \subseteq (\mathbb{N} \times \mathbb{D}^n)^*$ definable in STL, it holds that:*

- \mathcal{L} is safety iff there exists a formula ϕ of $G(\text{ppSTL})$ such that $\mathcal{L}(\phi) = \mathcal{L}$;
- \mathcal{L} is cosafety iff there exists a formula ϕ of $F(\text{ppSTL})$ such that $\mathcal{L}(\phi) = \mathcal{L}$.

Proof. We prove the theorem for the safety case (the proof for the cosafety fragment is identical). Let $\mathcal{L} \subseteq (\mathbb{N} \times \mathbb{D}^n)^\omega$ be a safety multi-valued language and let $\psi \in \text{STL}$ such that $\mathcal{L}(\psi) = \mathcal{L}$. *Ad absurdum*, we suppose that, for all $\phi \in G(\text{ppSTL})$, it holds that $\mathcal{L}(\phi) \neq \mathcal{L}(\psi)$.

We consider the formula $\mathbb{B}(\psi)$. By Definition 8, $\mathbb{B}(\psi)$ is a MTL formula over some set of atomic propositions \mathcal{AP} and $\mathcal{L}(\mathbb{B}(\psi)) \subseteq (\mathbb{N} \times 2^{\mathcal{AP}})^n$. Moreover, since $\mathcal{L}(\psi)$ is safety, by Lemma 2, it holds that $\mathcal{L}(\mathbb{B}(\psi))$ is a safety language as well. By the completeness theorem of the safety fragment of LTL [Chang *et al.*, 1992], there exists a formula $\psi' \in G(\text{ppMTL})$ such that $\mathcal{L}(\psi') = \mathcal{L}(\mathbb{B}(\psi))$. We remark that the theorem in [Chang *et al.*, 1992] is written for the case of LTL, but we can safely use it for MTL over qualitative time since the only thing that changes between LTL and MTL over qualitative time are the bounds on the operators, that do not alter the expressive power (but only the succinctness of the formalisms, that it is not of interest here).

Now, let ψ'' be the STL formula $\mathbb{B}^{-1}(\psi')$. It is easy to see that, since ψ' is of the form $G(\alpha)$ for some $\alpha \in \text{ppMTL}$, formula $\mathbb{B}^{-1}(\psi')$ is of the form $G(\alpha)$ for some $\alpha \in \text{ppSTL}$. Therefore, $\mathbb{B}^{-1}(\psi') \in G(\text{ppSTL})$. It follows that $\mathcal{L}(\psi'') = \mathcal{L}(\phi)$ for some formula ϕ in $G(\text{ppSTL})$. This is a contradiction with our hypothesis. Therefore, it has to hold that there exists a formula ϕ in $G(\text{ppSTL})$ such that $\mathcal{L}(\phi) = \mathcal{L}$. \square

B More examples for $G(\text{ppSTL})$ and $F(\text{ppSTL})$

Suppose to have a set of grants g_1, \dots, g_n and consider the following requirement: “if the arbiter is initially in configuration 1, then no grant can be issued; if, otherwise, the arbiter is initially in configuration 2, then there cannot be simultaneously issued two or more grants (mutual exclusion)”. Supposing that there is a variable x_c such that $x_c > 0$ if the configuration is 1 and $x_c \leq 0$ if the configuration is 2, the $G(\text{ppSTL})$ formula for the requirement is the following:

$$G(O(\tilde{Y} \perp \wedge x_c > 0) \rightarrow H(\bigwedge_{i=1}^n x_{g_i} \leq 0)) \wedge$$

$$G(O(\tilde{Y} \perp \wedge x_c \leq 0) \rightarrow H(\bigwedge_{i=1}^n \bigwedge_{j=1 \neq i}^n (x_{g_i} < 0 \vee x_{g_j} < 0)))$$

Note the crucial use of the subformula $\tilde{Y} \perp$ to hook the initial state of a trace. This requirement is *not definable* in the bounded fragment of STL.

C Considered datasets

We rely on the datasets Backblaze Hard Drive⁴, Tennessee Eastman Process⁵, and NASA C-MAPSS⁶.

⁴<https://www.backblaze.com/b2/hard-drive-test-data.html>

⁵<https://doi.org/10.7910/DVN/6C3JR1>

⁶<https://data.nasa.gov/dataset/C-MAPSS-Aircraft-Engine-Simulator-Data/xaut-bemq>

The *Backblaze Hard Drive* dataset contains continuously updated information on the “health” status of hard drives in the Backblaze data center, tracked employing Self Monitoring Analysis and Reporting Technology (SMART). Each trace is described by: date of the report, serial number of the drive, a label indicating a drive failure and 21 numerical SMART parameters. For comparison with the literature we focus on the ST4000DM000 hard drive model, and on a training set ranging from October to November 2016 (11423 good drives, 148 failing drives), and a test set composed of December 2016 recordings (23313 good drives, 86 failing drives), i.e., [Brunello *et al.*, 2023] Split *S1*. Note that, for time constraints, here we disregard the dataset [Brunello *et al.*, 2023] Split *S2*, which is a larger, although much less challenging version of Split *S1* (failure detection literature performance: F1-score ~ 0.93 vs F1-score ~ 0.56 , respectively).

The *Tennessee Eastman Process* (TEP) dataset is composed of simulated data from a fictitious chemical plant. It includes 1000 training and 1000 test traces, sampled every 3 minutes, and labelled with Type 0 (normal behaviour) or Type 1 (faulty behaviour). Each training (resp., test) trace lasts for 25 (resp., 48) hours. There are 500 faulty traces in both sets. Each trace has the features: trace ID, fault type, and 52 variables tracking data about the operating values of plant components.

The *NASA Commercial Modular Aero-Propulsion System Simulation* (C-MAPSS) dataset includes run-to-failure simulated data of turbofan jet engines. Specifically, in our considered dataset FD001, engines are simulated according to a single operating condition (called *Sea level*) and their failures are attributable to one possible cause (HPC degradation). Each engine simulation is represented by a multivariate time series, sampled at one value per second, obtained from 21 engine sensors. The dataset includes 100 training traces, each ending with a failure, and 100 test traces, each ending an arbitrary and known number of time steps before the failure (gap). To compare our framework with the literature [Kim and Sohn, 2020], and, specifically, with the detection performance of the *Unhealthy state* class, we proceed as follows. In the training set, 200 traces are generated from the original ones by considering the 70% prefix of a trace as normal behavior, and the full trace as faulty. As for the test set, each of the 100 traces is labeled as good (if trace length $< 0.7 \cdot (\text{length} + \text{gap})$, i.e., 61 traces) or not (if trace length $\geq 0.7 \cdot (\text{length} + \text{gap})$, i.e., 39 traces).

D Framework hyperparameters and tuning

Table 2 recaps the hyperparameters of the framework, along with a synthetic description. Their tuning ranges and final values for each dataset are instead reported in Table 3.

As mentioned in the main paper, we performed, for each dataset, a tuning phase with the GP-bfSTL version of the framework. To do that, we randomly partitioned each training set into 90% training and 10% validation, according to a stratified approach based on the failure label and grouping instances so as not to fragment system traces between the two partitions.

The actual tuning phase was carried out in two subsequent

steps. First, we employed a coarse evaluation to eliminate clearly sub-optimal parameter values. This led us to discard $pop_{size} = 50$, $max_{far} = 0.1$, $patience \in \{10, 30\}$, and $r_{interval} = 50$. In addition, briefly experimenting with the max_{far} parameter made clear to us that the latter is strongly tied to the number of good traces available for training the framework. Indeed, in the dataset *Backblaze*, which contains 11423 good traces, even setting $max_{far} = 0.02$ would allow the genetic algorithm to return a formula that had issued a bad prediction on over 200 good instances, which is disproportionate given the presence of only 148 failure traces. Thus, considering the most balanced situation of dataset C-MAPSS (100 good and 100 bad traces) as our reference scenario, we scaled the search range as follows: given i a max_{far} value, we consider $i/(max_{far}/100)$. A similar reasoning is applied to the parameter $fract_{good}$. Subsequently, we narrowed our focus to the refined search space. We relied on the Hyperopt-sklearn library⁷ with a budget of 50 iterations, since performing a grid search over all the remaining 2592 possibilities would have been still unfeasible. The performance of each tested parameter combination was the result of the average of 3 different framework executions. The monitored parameter was the F1-score, to be maximized.

E Extended experimental evaluations

In the following, we report an overview of the metrics we considered as well as some in-depth results regarding the experimental validation.

E.1 Metrics description

We employed an extensive array of metrics, that provide distinct insights into the framework’s failure detection performance: Precision, Recall, False Alarm Rate (FAR), F1 Score, and Matthews Correlation Coefficient (MCC). They are derived from the elements of a confusion matrix, where positive (resp., negative) instances correspond to failures (resp., good) traces: True Positives (TP), False Positives (FP), True Negatives (TN), and False Negatives (FN).

Precision and Recall are critical in classification tasks, like the one we considered. Precision indicates the proportion of identified failures that were actual failures. It is defined as:

$$\text{Precision} = \frac{TP}{TP + FP}.$$

Recall reflects the proportion of actual failures correctly identified, and it is essential in scenarios where missing a failure can have serious implications. It is defined as:

$$\text{Recall} = \frac{TP}{TP + FN}.$$

The F1 Score is the harmonic mean of Precision and Recall. This metric is valuable when balancing the need to identify as many failures as possible (high recall) against ensuring the identified cases are indeed failures (high precision). It is calculated as:

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}.$$

FAR is a critical metric in scenarios where the cost of false alarms is high. It indicates the likelihood of falsely identifying a non-failure as a failure. It is defined as:

$$\text{FAR} = \frac{FP}{FP + TN}.$$

As mentioned in the main paper, a low FAR is essential to maintain the reliability of the framework’s pool of formulas. Yet, note that FAR depends on TN, making it potentially misleading in scenarios where there is a high number of good traces and only a few, possibly hard to detect, failing ones.

Finally, MCC is defined as:

$$\text{MCC} = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}.$$

It ranges from -1 to +1, where +1 indicates perfect prediction, 0 no better than random prediction, and -1 total disagreement between prediction and actual value. MCC considers all four categories of the confusion matrix, thus providing a balanced metric to measure performance.

The choice of including the MCC in our analyses is motivated by its efficacy in handling class imbalance, as in the dataset *Backblaze S1*, offering a comprehensive view of the framework’s predictive performance [Chicco and Jurman, 2020]. Instead, the scarcity of failure events can lead to misleading results in the other considered metrics.

E.2 Additional results

The results of both our GP-bfSTL and GP-G(ppSTL) experiments are from 5 different runs. Of course, the same random seeds were employed for the two versions of the framework. Seeding is necessary, since the execution of the genetic programming algorithm is inherently stochastic. In addition, also the order in which failure traces are presented to the framework during its training phase (Algorithm 1, Line 4) is randomly determined, and may thus influence the number and kind of properties extracted.

In the following, we report some in-depth results regarding the GP-G(ppSTL) version of the framework, being the one based on the newly isolated well-founded monitorable logical fragment and also the one that achieved the overall best performance on the considered datasets.

Figure 5 shows, for each run, the evolution of the pool size throughout the training iterations. As can be seen, the pool size increases over time and eventually starts to stabilize. This stabilization likely occurs once no more useful predictive information can be extracted from the dataset. Still, the caveat already discussed in the main paper should be kept into account: better results may be achievable by performing multiple passes on the training set.

Finally, let us consider how the number of teacher forcing activations evolves as the training progresses. Intuitively, as the pool of properties expands, the frequency of these interventions should decrease. To confirm that, we proceed as

⁷<https://github.com/hyperopt/hyperopt-sklearn>

Table 2: Framework hyper parameters

Kind	Name	Description
General	n_{aug}	how many perturbed traces to generate for each trace, used both for \mathcal{T} and Σ_{\perp}
	min_{acc}	quality requirement for formula extraction: minimum accuracy over sub-traces of \mathcal{T}
	max_{far}	quality requirement for formula extraction: maximum FAR over traces of Σ_{\perp}
Genetic algorithm	$fract_{good}$	relative size of sets $\Sigma'_{\perp}, \Sigma''_{\perp}$ with respect to the size of Σ_{\perp}
	$r_{interval}$	rotation frequency, in number of generations, of Σ'_{\perp}
	max_{gen}	maximum number of generations
	$patience$	patience, in number of generations, of the early stopping criterion
	pop_{size}	number of individuals in the population
	$mutation_{prob}$	mutation probability, which undergoes the following decay over the generations: $mutation_{prob} / \sqrt[3]{generation_index}$
	$crossover_{prob}$	crossover probability

Table 3: Hyper parameters tuning ranges and final values

Kind	Name	Search range	Backblaze	C-MAPSS	TEP
General	n_{aug}	[50]	50	50	50
	min_{acc}	[0.5, 0.75]	0.75	0.75	0.75
	max_{far}	[0.0, 0.02, 0.05, 0.1]*	0.05*	0.02*	0.0*
Genetic algorithm	$fract_{good}$	[0.33, 0.5, 1.0, 2.0]*	2.0*	0.33*	0.5*
	$r_{interval}$	[1, 10, 50]	1	10	10
	max_{gen}	[500]	500	500	500
	$patience$	[10, 30, 50]	50	50	50
	pop_{size}	[50, 100, 200, 500]	200	500	200
	$mutation_{prob}$	[0.4, 0.5, 0.6, 0.7]	0.50	0.60	0.60
	$crossover_{prob}$	[0.6, 0.7, 0.8, 0.9]	0.70	0.80	0.80

*: raw max_{far} and $fract_{good}$ values, to be adjusted by the number of good training instances (num_{good}) present in a dataset according to the formula $param / (num_{good}/100)$, where $param \in \{max_{far}, fract_{good}\}$.

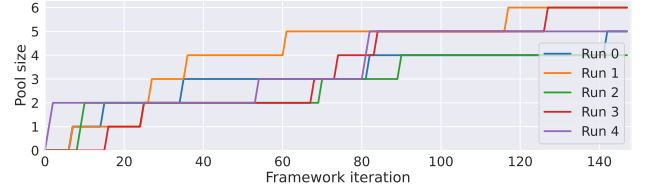
follows. For each run of the framework over a given dataset, we partition the iterations into intervals representing the 10% of the total iterations. Then, for each interval, we compute the total number of activations. For example, with the 100 iterations of the C-MAPSS dataset, this approach results in 10 aggregated values, each made by the sum of 10 elements. Figure 6 presents the results (transparent lines). For simplicity, within the same figure, we also report the average of the 5 curves (black line). For the TEP and C-MAPSS datasets, teacher forcing interventions diminish rapidly, reaching zero. This means that the related pools of formulas adequately capture all negative behaviours in the training set. Still, the test set performance (cfr. Table 1) for C-MAPSS was inferior than TEP, indicating a possible greater discrepancy between its training and test splits, or the presence of an overfitting phenomenon. In contrast, for the Backblaze dataset, while teacher forcing interventions decrease over time, they remain significantly above zero. This should not come as a surprise, given the notably lower performance of the framework observed for this dataset.

E.3 Examples of formulas extracted

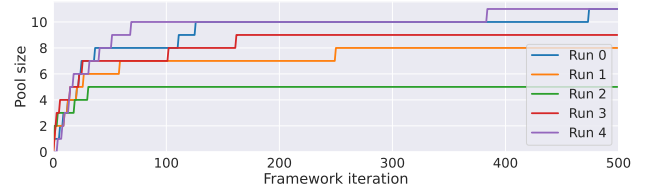
Here, we report some examples of formulas extracted by the framework during its training phase. These are not only highly interpretable but also effective in modeling the build-up, over time, of failure scenarios.

Backblaze S1

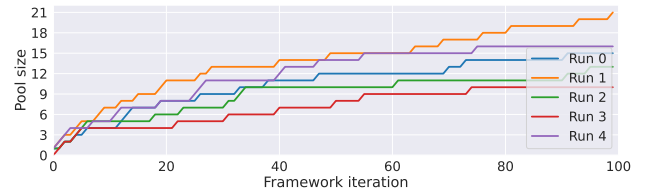
The monitoring pool \mathcal{P} includes the formula $G(\neg((smart_{198} \geq 1848.04) \wedge (smart_{198} \geq 30.86)))$, which computation tree is also described in Fig. 3. When it evaluates to false, a failure detection is triggered. The



(a) Backblaze S1 dataset.



(b) TEP dataset.



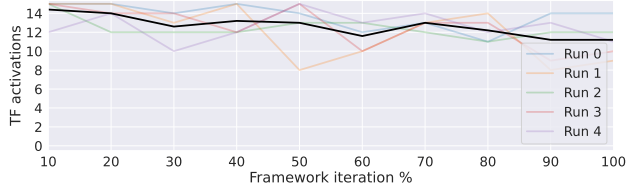
(c) C-MAPSS dataset.

Figure 5: Pool size increment along the framework training iterations, for 5 different runs of the framework (ppSTL formulas).

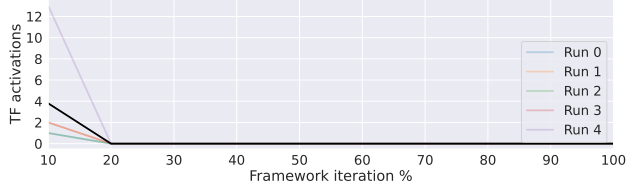
latter corresponds to a scenario where a sudden increase in the *uncorrectable sector count* (sensor $smart_{198}$) is witnessed. Indeed, following the SMART documentation, sensor $smart_{198}$ is particularly critical, as a rise in its value is likely to indicate defects of the disk surface and/or problems in the mechanical subsystem.

TEP

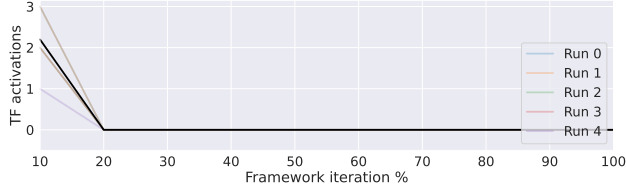
For the dataset TEP, we report a case where, first, the formula $f_1 = G(xmeas_1 < 0.494712 \vee xmeas_9 < 120.420)$ is extracted, encoding that a failure is going to happen if at any moment the *feed A (stream 1)* is greater than 0.494712 kscmh (kilo standard cubic meters per hour) while the *reac-*



(a) Backblaze S1 dataset.



(b) TEP dataset.



(c) C-MAPSS dataset.

Figure 6: Number of teacher forcing activations, aggregated every 10% training iterations, for 5 different runs of the framework (ppSTL formulas). Black line represents the average along the runs.

tor temperature exceeds 120.420°C . Next, exploiting f_1 , a new anticipatory formula $f_2 = G(H \neg ((xmv_2 \geq 53.9564) \wedge (xmeas_{29} \geq 34.5824)))$ is acquired. It signals a failure if there exists a time point in which the value of *component A of the purge gas analysis (stream 9)* is greater than $34.5824 \text{ kg h}^{-1}$ followed by an interval in which the *feed flow E (stream 3)* is steadily greater than 53.9564 mol\% .

C-MAPSS

Case 1 Formula $f_1 = G(O_{[0,1]}(sensor_3 < 1601.37 \wedge sensor_{14} < 8240.95))$ requires that, not to have a failure behaviour, at least once every two consecutive time instants, the *total temperature at HPC outlet* should be lower than 1601.37°R and the *corrected core speed* should be lower than 8240.95 rpm .

Case 2 In this scenario, first the formula $f_1 = G(H(sensor_{14} < 8143.33 \vee sensor_2 < 643.927))$ is extracted, requiring that, at any time, the *corrected core speed* must be lower than 8143.33 rpm or the *total temperature at LPC outlet* must be lower than 643.927°R . Such a formula, when violated, leads to the extraction of the formula $f_2 = G((sensor_2 < 643.573 \vee sensor_{11} < 48.1424))$, which requires either a tighter bound for *total temperature at LPC outlet* (it must be lower than 643.573°R) or that the *static pressure at HPC outlet* is lower than 48.1424 psia . Afterwards, based on the violation of f_2 , formula $f_3 =$

$G(sensor_{11} < 48.1188)$ is extracted, refining the threshold for the *static pressure at HPC outlet* to 48.1188 psia . Finally, upon violation of f_3 , formula $f_4 = G(O_{[0,3]}(sensor_{10} < 1.3 \wedge sensor_{15} < 8.49591))$ is learned. It requires that, at least once every four time instants, the *engine pressure ratio* must be lower than 1.3 while the *bypass ratio* does not exceed 8.49591 . It is worth highlighting how the iterative refinement of the properties, in this case, starts with an average preemptiveness of 0.45 seconds based on f_1 , increasing to 6.3 seconds after f_4 is learned.

F Reproducibility and code availability

We make available all the code to let interested people reproduce our experiments as well as use our framework in their domains for novel applications.

LINK: <https://github.com/dslab-uniud/ppSTL-IJCAI2024>

The material made available includes the code to learn a pool of formulas either using bfSTL or ppSTL, the dictionaries with the hyperparameters resulting from the tuning phase, the three datasets we employed, the Python packages (and their version) needed to reproduce the experiments, and the code to test the result of a framework execution. For detailed information on how to execute the code and/or reproduce the experiments, we refer the reader to the `readme.md` file contained in the home folder.