

Learning Interpretable Temporal Properties from Positive Examples Only

Rajarshi Roy¹, Jean-Raphaël Gaglione², Nasim Baharisangari³, Daniel Neider⁴, Zhe Xu³, and Ufuk Topcu²

¹ Max Planck Institute for Software Systems, Kaiserslautern, Germany

² University of Texas at Austin, Texas, USA

³ Arizona State University, Arizona, USA

⁴ Carl von Ossietzky University of Oldenburg, Oldenburg, Germany

Abstract. We consider the problem of explaining the temporal behavior of black-box systems using human-interpretable models. To this end, based on recent research trends, we rely on the fundamental yet interpretable models of *deterministic finite automata* (DFAs) and *linear temporal logic* (LTL) formulas. In contrast to most existing works for learning DFAs and LTL formulas, we rely on only positive examples. Our motivation is that negative examples are generally difficult to observe, in particular, from black-box systems. To learn meaningful models from positive examples only, we design algorithms that rely on *conciseness* and *language minimality* of models as regularizers. To this end, our algorithms adopt two approaches: a symbolic and a counterexample-guided one. While the symbolic approach exploits an efficient encoding of language minimality as a constraint satisfaction problem, the counterexample-guided one relies on generating suitable negative examples to prune the search. Both the approaches provide us with effective algorithms with theoretical guarantees on the learned models. To assess the effectiveness of our algorithms, we evaluate all of them on synthetic data.

Keywords: One-class learning · Automata learning · Learning of logic formulas.

1 Introduction

The recent surge of complex black-box systems in Artificial Intelligence has increased the demand for designing simple explanations of systems for human understanding. In several areas such as robotics, healthcare, and transportation (Royal-Society, 2019; Gunning et al., 2019; Molnar, 2022), inferring human-interpretable models has become the primary focus to promote human trust in systems.

To enhance the interpretability of systems, we aim at explaining their temporal behavior. To this end, models that are typically employed include, among others, finite state machines and temporal logics (Weiss, Goldberg, and Yahav, 2018; Roy, Fisman, and Neider, 2020). We focus here on the fundamental state machine: deterministic finite automata (DFAs) (Rabin and Scott, 1959); and formulas in the de facto standard temporal logic: linear temporal logic (LTL) (Pnueli, 1977). These models not only possess a host of desirable theoretical properties, but also feature easy-to-grasp syntax and intuitive semantics. The latter properties make them particularly suitable as interpretable models with many applications, e.g., as task knowledge for robotic agents (Kasenberg and Scheutz, 2017), as a formal specification for verification (Lemieux, Park, and Beschastnikh, 2015), as behavior classifier for unseen data (Shvo et al., 2021), and many more (Camacho and McIlraith, 2019).

The area of learning DFAs and LTL formulas is well-studied with a plethora of existing works (see related works). Most of them tackle the typical binary classification problem (Gold, 1978) of learning concise DFAs or LTL formulas from a finite set of examples partitioned into a positive and a negative set. However, negative examples are hard to obtain in many scenarios. Often, in safety-critical areas, observing negative examples from systems (e.g., from medical devices and self-driving cars) can be unrealistic (e.g., by injuring patients or hitting pedestrians). Further, often one only has access to a black-box implementation of the system and thus, can extract only its possible (i.e., positive) executions.

In spite of being relevant, the problem of learning concise DFAs and LTL formulas from positive examples, i.e., the corresponding *one class classification* (OCC) problem, has garnered little attention. The primary reason, we believe, is that, like most OCC problems, this problem is an ill-posed one. Specifically, the most concise model that classifies the positive examples correctly is the trivial model that classifies all examples as positive. This corresponds to a single state DFA or, in LTL, the formula *true*. These models, unfortunately, convey no insights about the underlying system.

To ensure a well-defined problem, (Avellaneda and Petrenko, 2018), who study the OCC problem for DFAs, propose the use of the (accepted) *language* of a model (i.e., the set of allowed executions) as a regularizer. Searching for a model that has a minimal language, however, results in the one that classifies only the given examples as positive. To avoid this overfitting, they additionally impose an upper bound on the size of the model. Thus, the OCC problem that they state is the following: given a set of positive examples P and a size bound n , learn a DFA that (a) classifies P correctly, (b) has size at most n , and (c) is language minimal. For language comparison, the order chosen is set inclusion.

To solve this OCC problem, Avellaneda and Petrenko (2018) then propose a *counterexample-guided* algorithm. This algorithm relies on generating suitable negative examples (i.e., counterexamples) iteratively to guide the learning process. Since only the negative examples dictate the algorithm, in many iterations of their algorithm, the learned DFAs do not have a language smaller (in terms of inclusion) than the previous hypothesis DFAs. This results in searching through several unnecessary DFAs.

To alleviate this drawback, as our first contribution, we present a *symbolic* algorithm for solving the OCC problem for DFA. Our algorithm converts the search for a language minimal DFA symbolically to a series of satisfiability problems in Boolean propositional logic. Our novelty includes an efficient encoding of the language inclusion check of DFAs in a propositional formula, which is polynomial in the size of the DFAs. We then exploit an off-the-shelf SAT solver to check satisfiability of the generated propositional formulas and subsequently, construct a suitable DFA. We expand on this algorithm in Section 3.

We then present two novel algorithms for solving the OCC problem for LTL formulas. Our first algorithm is a *semi-symbolic* algorithm, which combines ideas from both the symbolic and the counterexample-guided approaches. Roughly, this algorithm exploits negative examples to overcome the theoretical difficulties of symbolically encoding language inclusion for LTL (since LTL inclusion check is known to be inherently harder than that for DFAs (Sistla and Clarke, 1985)). Our second algorithm is simply a counterexample-guided algorithm that relies solely on the generation of negative examples for learning. We detail on both the algorithms in Section 4.

Finally, we empirically evaluate all of our algorithms on synthetic benchmarks to assess their efficacy. We demonstrate that our symbolic algorithm solves the OCC problem for DFA in fewer iterations and runtime comparable to the counterexample-guided algorithm, skipping thousands of counterexample generation. Our results show that the average ratio of the number of iterations of our symbolic algorithm to the counterexample-guided algorithm is 0.14 while the average ratio of the inference time of our symbolic algorithm to the counterexample-guided algorithm is 1.09. We also demonstrate that our symbolic algorithm solves the OCC problem for LTL faster than the counterexample-guided algorithm does. Our results show that our semi-symbolic algorithm is at most 260.73% faster than the counterexample-guided algorithm for inferring LTL formulas.

Related Works. The OCC problem described in this paper belongs to the body of works categorized as passive learning (Gold, 1978). As alluded to in the introduction, in this topic, the most popular problem is the binary classification problem for learning DFAs and LTL formulas. Notable works include the works by Biermann and Feldman (1972); Grinchtein, Leucker, and Piterman (2006); Heule and Verwer (2010) for DFAs, and Neider and Gavran (2018); Camacho and McIlraith (2019); Raha et al. (2022) for LTL formulas.

The OCC problem of learning formal models from positive examples was first studied by Gold (1967). This work showed that the exact identification of certain models (that includes DFAs and LTL formulas) from positive examples is not possible. Thereby, works have mostly focussed on models that are learnable easily from positive examples, such as pattern languages (Angluin, 1980), stochastic finite state machines (Carrasco and Oncina, 1999), and hidden Markov Models (Stolcke and Omohundro, 1992). None of these works considered learning DFAs or LTL formulas, mainly due to the lack of a meaningful regularizer.

Recently, Avellaneda and Petrenko (2018) proposed the use of language minimality as a regularizer and thereafter, developed an effective algorithm for learning DFAs. While their algorithm cannot overcome the theoretical difficulties shown by Gold (1967), they still produce a DFA that is a concise description of the positive examples. We significantly improve upon their algorithm by relying on a novel encoding of language minimality using propositional logic.

For temporal logics, there are a few works that consider the OCC problem. Notably, Ehlers, Gavran, and Neider (2020) propose a learning algorithm for a fragment of LTL which permits an automata representation known as universally very-weak automata (UVWs). However, since their algorithm relies on UVWs, which has strictly less ex-

pressive power than LTL, it cannot be extended to full LTL. Further, there are works on learning LTL (Chou, Ozay, and Berenson, 2022) and STL (Jha et al., 2019) formulas from trajectories of high-dimensional systems. These works base their learning on the assumption that the underlying system optimizes some cost function. Our method, in contrast, is based on the natural notion of language minimality to find tight descriptions, without any assumptions about the system.

A problem similar to our OCC problem is studied in the context of inverse reinforcement learning (IRL) to learn temporal rewards for RL agents from (positive) demonstrations. For example, Kasenberg and Scheutz (2017) learn concise LTL formulas based on a probabilistic measure that captures how “well” a formula distinguishes the provided demonstrations from random executions of the system. To generate the random executions, they rely on a Markov Decision Process (MDP) implementation of the underlying system. Our regularizers, in contrast, assume the underlying system to be a black-box and need no access to its internal mechanisms. Vazquez-Chanlatte et al. (2018) also learn LTL-like formulas from demonstrations. Their search requires a pre-computation of the lattice of formulas induced by the subset order, which can be a bottleneck for scaling to full LTL.

2 Preliminaries

In this section, we set up the notation for the rest of the paper.

Let $\mathbb{N} = \{1, 2, \dots\}$ be the set of natural numbers and $[n] = \{1, 2, \dots, n\}$ be the set of natural numbers up to n .

Words and Languages. To formally represent system executions, we rely on the notion of words, defined over a finite and nonempty *alphabet* Σ . The elements of Σ , which denote relevant system states, are referred to as *symbols*.

A *word* over Σ is a finite sequence $w = a_1 \dots a_n$ with $a_i \in \Sigma$, $i \in [n]$. The *empty word* ε is the empty sequence. The length $|w|$ of w is the number of symbols in it (note that $|\varepsilon| = 0$). Moreover, Σ^* denotes the set of all words over Σ . Further, we use $w[i] = a_i$ to denote the i -th symbol of w and $w[i:] = a_i \dots a_n$ to denote the suffix of w starting from position i .

A *language* L is any set of words from Σ^* . We allow the standard set operations on languages such as $L_1 \subseteq L_2$, $L_1 \subset L_2$ and $L_1 \setminus L_2$.

Propositional logic. All our algorithms rely on propositional logic and thus, we introduce it briefly. Let Var be a set of propositional variables, which take Boolean values $\{0, 1\}$ (0 represents *false*, 1 represents *true*). Formulas in propositional logic—which we denote by capital Greek letters—are defined recursively as

$$\Phi := x \in Var \mid \neg\Phi \mid \Phi \vee \Psi$$

. As syntax sugar, we allow the following standard formulas: *true*, *false*, $\Phi \wedge \Psi$, $\Phi \rightarrow \Psi$ and $\Phi \leftrightarrow \Psi$.

An *assignment* $v: Var \mapsto \{0, 1\}$ maps propositional variables to Boolean values. Based on an assignment v , we define the semantics of propositional logic using a valuation function $V(v, \Phi)$, which is inductively defined as follows:

$$\begin{aligned} V(v, x) &= v(x) \\ V(v, \neg\Psi) &= 1 - V(v, \Psi) \\ V(v, \Psi \vee \Phi) &= \max\{V(v, \Psi), V(v, \Phi)\} \end{aligned}$$

We say that v satisfies Φ if $V(v, \Phi) = 1$, and call v a model of Φ . A formula Φ is satisfiable if there exists a model v of Φ .

Arguably, the most well-known problem in propositional logic—the satisfiability (SAT) problem—is the problem of determining whether a propositional formula is satisfiable or not. With the rapid development of SAT solvers (Li and Manyà, 2021), checking satisfiability of formulas with even millions of variables is feasible. Most solvers can also return a model when a formula is satisfiable.

3 Learning DFA from Positive Examples

In this section, we present our symbolic algorithm for learning DFAs from positive examples. We begin by formalizing DFAs.

A *deterministic finite automaton* (DFA) is a tuple $\mathcal{A} = (Q, \Sigma, \delta, q_I, F)$ where Q is a finite set of states, Σ is the alphabet, $q_I \in Q$ is the initial state, $F \subseteq Q$ is the set of final states, and $\delta: Q \times \Sigma \rightarrow Q$ is the transition function. We define the size $|\mathcal{A}|$ of a DFA as its number of states $|Q|$.

Given a word $w = a_1 \dots a_n \in \Sigma^*$, the run of \mathcal{A} on w , denoted by $\mathcal{A}: q_1 \xrightarrow{w} q_{n+1}$, is a sequence of states and symbols $q_1 a_1 q_2 a_2 \dots a_n q_{n+1}$, such that $q_1 = q_I$ and for each $i \in [n]$, $q_{i+1} = \delta(q_i, a_i)$. Moreover, we say w is accepted by \mathcal{A} if $q_{n+1} \in F$. Finally, we define the language of \mathcal{A} as $L(\mathcal{A}) = \{w \in \Sigma^* \mid w \text{ is accepted by } \mathcal{A}\}$.

To introduce the OCC problem for DFAs, we first describe the learning setting. The OCC problem relies on a set of positive examples, which we represent using a finite set of words $P \subset \Sigma^*$. Additionally, the problem allows a bound n to restrict the size of the learned DFA. The role of this size bound is two-fold. First, it ensures that the learned DFA does not overfit P . Second, using a suitable bound, one can force the learned DFAs to be concise (and thus, interpretable).

Finally, we define a DFA \mathcal{A} to be an n -description of P if $P \subseteq L(\mathcal{A})$ and $|\mathcal{A}| \leq n$. When clear from the context, we simply say n -description, instead of n -description of P .

Having set up the prerequisites, we now state the OCC problem for DFA as follows:

Problem 1 (OCC problem for DFAs). Given a set of positive words P and a size bound n , learn a DFA \mathcal{A} such that:

- (1) \mathcal{A} is an n -description; and
- (2) for all DFA \mathcal{A}' that is an n -description, $L(\mathcal{A}') \not\subseteq L(\mathcal{A})$.

Intuitively, the problem stated above asks to search for an n -description DFA with a minimal language. Note that several such DFAs can exist since the language inclusion order used to determine minimality is only a partial order on the languages of DFA. We, here, are interested in learning only one such DFA, leaving the problem of learning all such DFAs as interesting future work.

3.1 The Symbolic Algorithm

We now present our algorithm for solving Problem 1. The underlying idea of our algorithm is to reduce the search for an appropriate DFA to a series of satisfiability checks of propositional formulas. Each satisfiable propositional formula helps us to construct a hypothesis DFA \mathcal{A} , which is the current guess. Based on the hypothesis \mathcal{A} , we construct a propositional formula $\Phi^{\mathcal{A}}$ which enables the search for the next hypothesis \mathcal{A}' that has a language smaller (in the inclusion order) than the current one. Precisely, we construct a propositional formula $\Phi^{\mathcal{A}}$ that has the following properties: 1. $\Phi^{\mathcal{A}}$ is satisfiable if and only if there exists a DFA \mathcal{A}' that is an n -description and $L(\mathcal{A}') \subset L(\mathcal{A})$; and 2. based on a model v of $\Phi^{\mathcal{A}}$, one can construct such a DFA.

Based on this main ingredient $\Phi^{\mathcal{A}}$, we design our learning algorithm, which we sketch in Algorithm 1. Our algorithm initializes the hypothesis DFA \mathcal{A} to be \mathcal{A}_{Σ^*} , which is the one-state DFA that accepts all words in Σ^* . Observe that \mathcal{A}_{Σ^*} is trivially an n -description, since $P \subset \Sigma^*$ and $|\mathcal{A}_{\Sigma^*}| = 1$. The algorithm then iteratively exploits $\Phi^{\mathcal{A}}$ to construct the next hypothesis DFAs, until $\Phi^{\mathcal{A}}$ is unsatisfiable for some \mathcal{A} . On being unsatisfiable, we terminate and return the current hypothesis \mathcal{A} as a solution. The correctness of this algorithm follows from the following theorem:

Theorem 1. *Given positive words P and a size bound n , Algorithm 1 learns a DFA \mathcal{A} that is an n -description and for all DFA \mathcal{A}' that is an n -description, $L(\mathcal{A}') \not\subseteq L(\mathcal{A})$.*

We now turn to the construction of $\Phi^{\mathcal{A}}$. To achieve the aforementioned properties, we define $\Phi^{\mathcal{A}}$ as follows:

$$\Phi^{\mathcal{A}} := \Phi_{\text{DFA}} \wedge \Phi_P \wedge \Phi_{\subseteq \mathcal{A}} \wedge \Phi_{\not\subseteq \mathcal{A}} \quad (1)$$

The first conjunct Φ_{DFA} ensures that we learn a valid DFA \mathcal{A}' . The second conjunct Φ_P ensures that \mathcal{A}' accepts all positive words. The third conjunct $\Phi_{\subseteq \mathcal{A}}$ ensures that $L(\mathcal{A}')$ is a subset of $L(\mathcal{A})$. The final conjunct $\Phi_{\not\subseteq \mathcal{A}}$ ensures that

Algorithm 1: Symbolic Algorithm for Learning DFA

Input: Positive words P , bound n

- 1: $\mathcal{A} \leftarrow \mathcal{A}_{\Sigma^*}$, $\Phi^{\mathcal{A}} := \Phi_{\text{DFA}} \wedge \Phi_P$
 - 2: **while** $\Phi^{\mathcal{A}}$ is satisfiable (with model v) **do**
 - 3: $\mathcal{A} \leftarrow$ DFA constructed from v
 - 4: $\Phi^{\mathcal{A}} := \Phi_{\text{DFA}} \wedge \Phi_P \wedge \Phi_{\subseteq \mathcal{A}} \wedge \Phi_{\supseteq \mathcal{A}}$
 - 5: **end while**
 - 6: **return** \mathcal{A}
-

$L(\mathcal{A}')$ is not a superset of $L(\mathcal{A})$. Together, conjuncts $\Phi_{\subseteq \mathcal{A}}$ and $\Phi_{\supseteq \mathcal{A}}$ ensure that $L(\mathcal{A}')$ is a proper subset of $L(\mathcal{A})$. In what follows, we detail the construction of each conjunct.

To encode the hypothesis DFA $\mathcal{A}' = (Q', \Sigma, \delta', q_I', F')$ symbolically, we follow Heule and Verwer (2010) and rely on the propositional variables: (1) $d_{p,a,q}$ where $p, q \in [n]$ and $a \in \Sigma$; and (2) f_q where $q \in [n]$. The variables $d_{p,a,q}$ and f_q encode the transition function δ' and the final states F' , respectively, of \mathcal{A}' . Mathematically speaking, if $d_{p,a,q}$ is set to true, then $\delta'(p, a) = q$ and if f_q is set to true, then $q \in F'$. Note that we denote the states Q' using set $[n]$ and the initial state q_I' using numeral 1.

Now, to ensure \mathcal{A}' has a deterministic transition function δ' , Φ_{DFA} asserts the following constraint:

$$\bigwedge_{p \in [n]} \bigwedge_{a \in \Sigma} \left[\bigvee_{q \in [n]} d_{p,a,q} \wedge \bigwedge_{q \neq q' \in [n]} [\neg d_{p,a,q} \vee \neg d_{p,a,q'}] \right].$$

Based on a model v of the variables $d_{p,a,q}$ and f_q , we can simply construct \mathcal{A}' . We set $\delta'(p, a)$ to be the unique state q for which $v(d_{p,a,q}) = 1$ and $q \in F'$ if $v(f_q) = 1$.

Next, to construct conjunct Φ_P , we introduce variables $x_{u,q}$ where $u \in \text{Pref}(P)$ and $q \in [n]$, which track the run of \mathcal{A}' on all words in $\text{Pref}(P)$, which is the set of prefixes of all words in P . Precisely, if $x_{u,q}$ is set to true, then there is a run of \mathcal{A}' on u ending in state q , i.e., $\mathcal{A}' : q_I' \xrightarrow{u} q$.

Using the introduced variables, Φ_P ensures that the words in P are accepted by imposing the following constraints:

$$\begin{aligned} & x_{\varepsilon,1} \wedge \bigwedge_{q \in \{2, \dots, n\}} \neg x_{\varepsilon,q} \\ & \bigwedge_{u \in \text{Pref}(P)} \bigwedge_{p, q \in [n]} \bigwedge_{a \in \Sigma} [x_{u,p} \wedge d_{p,a,q} \rightarrow x_{ua,q}] \\ & \bigwedge_{w \in P} \bigwedge_{q \in [n]} x_{w,q} \rightarrow f_q \end{aligned}$$

The first constraint above ensures that the runs start in the initial state q_I' (denoted using numeral 1) while the second constraint ensures that the runs adhere to the transition function. The third constraint ensures that the run of \mathcal{A}' on every $w \in P$ ends in a final state, ensuring their acceptance.

For the third conjunct $\Phi_{\subseteq \mathcal{A}}$, we must track the synchronized runs of the current hypothesis \mathcal{A} and the next hypothesis \mathcal{A}' to compare their runs on all words in Σ^* . Towards this, we introduce some more variables, $y_{q,q'}^{\mathcal{A}}$ where $q, q' \in [n]$. Precisely, $y_{q,q'}^{\mathcal{A}}$ is set to true, if there exists a word $w \in \Sigma^*$ such that there are runs $\mathcal{A} : q_I \xrightarrow{w} q$ and $\mathcal{A}' : q_I' \xrightarrow{w} q'$.

To ensure $L(\mathcal{A}') \subseteq L(\mathcal{A})$, $\Phi_{\subseteq \mathcal{A}}$ imposes the following constraints:

$$\begin{aligned} & y_{1,1}^{\mathcal{A}} \\ & \bigwedge_{q=\delta(p,a)} \bigwedge_{p', q' \in [n]} \bigwedge_{a \in \Sigma} \left[[y_{p,p'}^{\mathcal{A}} \wedge d_{p',a,q'}] \rightarrow y_{q,q'}^{\mathcal{A}} \right] \\ & \bigwedge_{p \notin F} \bigwedge_{p' \in [n]} \left[y_{p,p'}^{\mathcal{A}} \rightarrow \neg f_{p'} \right] \end{aligned}$$

The first constraint ensures that the synchronized runs of \mathcal{A} and \mathcal{A}' start in the respective initial states, while the second constraint ensures that they adhere to their respective transition functions. The third constraint ensures that if the synchronized run ends in a non-final state in \mathcal{A} , it must also end in a non-final state in \mathcal{A}' .

The construction of the final conjunct $\Phi_{\mathcal{D},\mathcal{A}}$ relies on the following result:

Lemma 1. *Let $\mathcal{A}, \mathcal{A}'$ be DFAs such that $|\mathcal{A}| = |\mathcal{A}'| = n$, $L(\mathcal{A}') \subset L(\mathcal{A})$, and $K = n^2$. Then there exists a word $w \in \Sigma^*$ of length $\leq K$, such that $w \in L(\mathcal{A})$ and $w \notin L(\mathcal{A}')$.*

The above result, essentially, establishes an upper bound to the length of the word that distinguishes DFAs \mathcal{A} and \mathcal{A}' .

Based on this result, we introduce variables $z_{i,q,q'}$ where $i \in [n^2]$ and $q, q' \in [n]$, to track the synchronized run of \mathcal{A} and \mathcal{A}' on a word of length at most $K = n^2$. Precisely, if $z_{i,q,q'}$ is set to true, then there exists a word w of length i , with the runs $\mathcal{A}: q_I \xrightarrow{w} q$ and $\mathcal{A}': q'_I \xrightarrow{w} q'$.

Now, $\Phi_{\mathcal{D},\mathcal{A}}$ imposes the following constraints:

$$\begin{aligned} & \bigwedge_{i \in [n^2]} \left[\bigvee_{q, q' \in [n]} z_{i,q,q'} \wedge \left[\bigwedge_{\substack{p \neq q \in [n] \\ p' \neq q' \in [n]}} \neg z_{i,p,p'} \vee \neg z_{i,q,q'} \right] \right] \\ & \bigwedge_{\substack{p, q \in [n] \\ p', q' \in [n]}} \left[\left[z_{i,p,p'} \wedge z_{i+1,q,q'} \right] \rightarrow \bigvee_{\substack{a \in \Sigma \text{ where} \\ q = \delta(p,a)}} d_{p',a,q'} \right] \\ & \bigvee_{i \in [n^2]} \bigvee_{\substack{q \in [n] \\ q' \in F}} \left[z_{i,q,q'} \wedge \neg f_{q'} \right] \end{aligned}$$

The first three constraints above ensure that words up to length n^2 have a valid synchronised run on the two DFAs \mathcal{A} and \mathcal{A}' . The final constraint ensures that there is a word of length $\leq n^2$ on which the synchronized run ends in a final state in \mathcal{A} but not in \mathcal{A}' .

4 Learning LTL formulas from Positive Examples

We now switch our focus to algorithms for learning LTL formulas from positive examples. We begin with a formal introduction to LTL.

Linear temporal logic (LTL) is a logic that reasons about temporal behavior of systems using temporal modalities. While traditionally LTL is built over propositional variables \mathcal{P} Pnueli (1977), to unify the notation with DFAs, we define LTL over an alphabet Σ . It is, however, not a restriction since an LTL formula over \mathcal{P} can always be translated to an LTL formula over $\Sigma = 2^{\mathcal{P}}$. Formally, LTL formulas—which we denote by small Greek letters—are defined inductively as:

$$\varphi := a \in \Sigma \mid \neg\varphi \mid \varphi \vee \psi \mid \mathbf{X}\varphi \mid \varphi \mathbf{U} \psi$$

As syntax sugar, along with additional constants and operators used in propositional logic, we allow the standard temporal operators \mathbf{F} (“finally”) and \mathbf{G} (“globally”). We define $\Lambda = \{\neg, \vee, \wedge, \rightarrow, \mathbf{X}, \mathbf{U}, \mathbf{F}, \mathbf{G}\} \cup \Sigma$ as the set of all operators (which, for simplicity, also includes symbols). We define the size $|\varphi|$ of φ as the number of its unique subformulas; e.g., size of $\varphi = (a \mathbf{U} \mathbf{X} b) \vee \mathbf{X} b$ is five, since its five distinct subformulas are $a, b, \mathbf{X} b, a \mathbf{U} \mathbf{X} b$, and $(a \mathbf{U} \mathbf{X} b) \vee \mathbf{X} b$.

To interpret LTL formulas over (finite) words, we follow the semantics proposed by Giacomo and Vardi (2013). Given a word w , we define recursively when a LTL formula holds at position i , i.e., $w, i \models \varphi$, as follows:

$$\begin{aligned} w, i & \models a \in \Sigma \text{ if and only if } a = w[i] \\ w, i & \models \neg\varphi \text{ if and only if } w, i \not\models \varphi \\ w, i & \models \mathbf{X}\varphi \text{ if and only if } i < |w| \text{ and } w, i+1 \models \varphi \\ w, i & \models \varphi \mathbf{U} \psi \text{ if and only if } w, j \models \psi \text{ for some} \\ & \quad i \leq j \leq |w| \text{ and } w, i' \models \varphi \text{ for all } i \leq i' < j \end{aligned}$$

We say w satisfies φ or, alternatively, φ holds on w if $w, 0 \models \varphi$, which, in short, is written as $w \models \varphi$.

The OCC problem for LTL formulas, similar to Problem 1, relies upon a set of positive words $P \subset \Sigma^*$ and a size upper bound n . Moreover, an LTL formula φ is an n -description of P if, for all $w \in P$, $w \models \varphi$, and $|\varphi| \leq n$. For brevity, we omit P from n -description when obvious.

We state the OCC problem for LTL formulas as follows:

Problem 2 (OCC problem for LTL formulas). Given a set of positive words P and a size bound n , learn an LTL formula φ such that:

- (1) φ is an n -description; and
- (2) for all LTL formula φ' that is an n -description, $\varphi' \not\rightarrow \varphi$ or $\varphi \rightarrow \varphi'$.

Intuitively, the above problem asks to search for an LTL formula φ that is an n -description and holds on a minimal set of words. Once again, like Problem 1, there can be several such LTL formulas, but we are interested in learning exactly one of them.

4.1 The Semi-Symbolic Algorithm

This algorithm, being *semi*-symbolic, does not entirely depend on the current hypothesis, an LTL formula φ here, as was the case in Algorithm 1. In addition, it relies on a set of negative examples N , accumulated during the algorithm. Thus, relying on both the current hypothesis φ and the negative examples N , we construct a propositional formula $\Psi^{\varphi, N}$ to search for the next hypothesis φ' . Concretely, $\Psi^{\varphi, N}$ has the properties that: (1) $\Psi^{\varphi, N}$ is satisfiable if and only if there exists an LTL formula φ' that is an n -description, does not hold on words in N , and $\varphi \not\rightarrow \varphi'$; and (2) based on a model v of $\Psi^{\varphi, N}$, one can construct such an LTL formula φ^v .

The semi-symbolic algorithm follows a paradigm similar to the one illustrated in Algorithm 1. However, unlike the previous algorithm, the current guess φ' , obtained from a model of $\Psi^{\varphi, N}$, may not always satisfy the relation $\varphi' \rightarrow \varphi$. In such a case, we generate a word (i.e., a negative example) that satisfies φ , but not φ' to eliminate φ' from the search space. The generation of words is performed by constructing DFAs from the LTL formulas (Zhu et al., 2017) and then performing a breadth-first search over them. If, otherwise, $\varphi' \rightarrow \varphi$, we then update our current hypothesis and continue until $\Psi^{\varphi, N}$ is unsatisfiable. Overall, this algorithm learns an appropriate LTL formula with the following guarantee:

Theorem 2. *Given positive words P and size bound n , the semi-symbolic algorithm learns an LTL formula φ that is an n -description and for all LTL formulas φ' that is an n -description, $\varphi' \not\rightarrow \varphi$ or $\varphi \rightarrow \varphi'$.*

We now focus on the construction of $\Psi^{\varphi, N}$, which is significantly different from that of Φ^A . It is defined as follows:

$$\Psi^{\varphi, N} := \Psi_{\text{LTL}} \wedge \Psi_P \wedge \Psi_N \wedge \Psi_{\neq \varphi}. \quad (2)$$

The first conjunct Ψ_{LTL} ensures that the hypothesis is a valid LTL formula φ' . The second conjunct Ψ_P ensures that φ' holds on all positive words, while the third, Ψ_N , ensures that it does not hold on the negative words. The final conjunct $\Psi_{\neq \varphi}$ ensures that $\varphi \not\rightarrow \varphi'$.

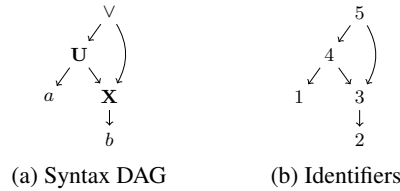


Fig. 1: Syntax DAG and identifiers for $(a \text{ U } X b) \vee X b$

Following Neider and Gavran (2018), all of our conjuncts rely on a canonical syntactic representation of LTL formulas as *syntax DAGs*. A syntax DAG is a directed acyclic graph (DAG) that is obtained from the syntax tree of an LTL formula by merging its common subformulas. An example of a syntax DAG is illustrated in Figure 1a. Further, to uniquely identify each node of a syntax DAG, we assign them unique identifiers from $[n]$ in a way that every parent node has an identifier larger than its children (see Figure 1b).

Now, to construct the hypothesis φ' , we encode its syntax DAG, using the following propositional variables: (1) $x_{i,\lambda}$ for $i \in [n]$ and $\lambda \in \Lambda$; and (2) $l_{i,j}$ and $r_{i,j}$ for $i \in [n]$ and $j \in [i-1]$. The variable $x_{i,\lambda}$ tracks the operator label of the Node i of the syntax DAG of φ' , while variables $l_{i,j}$ and $r_{i,j}$ track the left and right child, respectively, of Node i . Mathematically, $x_{i,\lambda}$ is set to true if and only if Node i is labeled with operator λ . Precisely, if $l_{i,j}$ (resp. $r_{i,j}$) is set to true if and only if left (resp. right) child of Node i is Node j .

To ensure variables $x_{i,\lambda}$, $l_{i,j}$ and $r_{i,j}$ have the desired meaning, Ψ_{LTL} imposes certain structural constraints. For instance, to ensure that each node of the syntax DAG of φ' is uniquely labeled by an operator, we have the following constraints:

$$\bigwedge_{i \in [n]} \left[\bigvee_{\lambda \in \Lambda} x_{i,\lambda} \wedge \bigwedge_{\lambda \neq \lambda' \in \Lambda} [\neg x_{i,\lambda} \vee \neg x_{i,\lambda'}] \right].$$

Ψ_{LTL} has more structural constraints for which we refer the readers to Neider and Gavran (2018).

We now move on to the construction of Ψ_P and Ψ_N . They relies on variables $y_{w,t}^i$ where $i \in [n]$, $w \in P \cup N$, and $t \in [|w|]$. The variable $y_{w,t}^i$ tracks whether $\varphi'[i]$ holds on suffix $w[t:]$, where $\varphi'[i]$ is the subformula of φ' rooted at Node i . Precisely, $y_{w,t}^i$ is set to true if and only if $w[t:] \models \varphi'[i]$.

To ensure desired meaning of variables $y_{w,t}^i$, Ψ_P and Ψ_N rely on semantic constraints, again similar to ones proposed by Neider and Gavran (2018). Exemplarily, the constraint implementing the semantics of the \mathbf{X} -operator is the following:

$$\bigwedge_{\substack{i \in [n] \\ j \in [i-1]}} x_{i,\mathbf{X}} \wedge l_{i,j} \rightarrow \left[\bigwedge_{t \in [|w|-1]} y_{w,t}^i \leftrightarrow y_{w,t+1}^j \right] \wedge \neg y_{w,|w|}^i.$$

Intuitively, this constraint states that if Node i is labeled with \mathbf{X} and Node j is the left child of Node i , then $\varphi'[i]$ holds on $w[t:]$ if and only if $\varphi'[j]$ holds on $w[t+1:]$ (i.e., if $t < |w|$). For the other LTL operators, we have similar semantic constraints, which the readers can find in the supplementary material.

Now to ensure that φ' holds on positive words Ψ_P additionally imposes $\bigwedge_{w \in P} y_{w,0}^n$, while, to ensure φ' does not hold on negative words, Ψ_N imposes $\bigwedge_{w \in N} \neg y_{w,0}^n$.

The conjunct $\Psi_{\neq \varphi}$ requires symbolically encoding a word u of length upto a *time horizon* of $K = 2^{2^{n+1}}$. The choice of K is derived from the fact that size of equivalent DFAs for an LTL formulas can be at most doubly exponential Giacomo and Vardi (2015) and the result in Lemma 1. Now our encoding of a symbolic word relies on variables $p_{t,a}$ where $t \in [K]$ and $a \in \Sigma$. If $p_{t,a}$ is set to true, then $u[t] = a$.

To ensure that the variables $p_{t,a}$ encode a valid word u , we generate a formula Ψ_{word} that has constraints like the following:

$$\bigwedge_{t \in [K]} \left[\bigvee_{a \in \Sigma \cup \{\epsilon\}} p_{t,a} \wedge \bigwedge_{a \neq a' \in \Sigma \cup \{\epsilon\}} [\neg p_{t,a} \vee \neg p_{t,a'}] \right].$$

The above constraint ensures that, in the word u , each position $t \leq K$ has a unique symbol from $\Sigma \cup \{\epsilon\}$.

Further, to track whether φ and φ' hold on u , we have variables $z_{u,t}^{\varphi,i}$ and $z_{u,t}^{\varphi',i}$ where $i \in [n]$, $t \in [K]$. These variables are similar to $y_{w,t}^i$, in the sense that $z_{u,t}^{\varphi,i}$ (resp. $z_{u,t}^{\varphi',i}$) is set to true, if φ (resp. φ') holds at position t . To assign meaning to these variables, we impose constraints, which we denote using Ψ_{sem} , similar to the semantic constraints imposed on $y_{w,t}^i$. For instance, for a symbol a , we have the following constraint:

$$\bigwedge_{i \in [n]} \bigwedge_{a \in \Sigma} \left[x_{i,a} \rightarrow \left[\bigwedge_{t \in [K]} z_{u,t}^{\varphi',i} \leftrightarrow p_{t,a} \right] \right].$$

Finally, we express $\Psi_{\neq\varphi}$ as follows:

$$\Psi_{\text{word}} \wedge \Psi_{\text{sem}} \wedge z_{u,0}^{\varphi,n} \wedge \neg z_{u,0}^{\varphi',n}$$

Intuitively, this constraint ensures that there exists a word on which φ holds and φ' does not.

4.2 The Counterexample-guided Algorithm

We now design a *counterexample-guided* algorithm to solve Problem 2. In contrast to the symbolic (or semi-symbolic) algorithm, this algorithm does not guide the search based on propositional formulas built out of the hypothesis LTL formula. Instead this algorithm relies entirely on two sets: a set of negative words N and a set of discarded LTL formulas D . Based on these two sets, we design a propositional formula $\Omega^{N,D}$ that has the properties that: (1) $\Omega^{N,D}$ is satisfiable if and only if there exists an LTL formula φ that is an n -description, does not hold on $w \in N$, and is not one of the formulas in D ; and (2) based on a model v of $\Omega^{N,D}$, one can construct such an LTL formula.

Being a counterexample-guided algorithm, the construction of the sets N and D forms the crux of the algorithm. In each iteration, these sets are updated based on the relation between the hypothesis φ and the current guess φ' (obtained from a model of $\Omega^{N,D}$). There are exactly three relevant cases, which we discuss briefly.

- First, $\varphi' \leftrightarrow \varphi$, i.e., φ' and φ hold on the exact same set of words. In this case, the algorithm discards φ' by adding it to D due to its similarity to φ .
- Second, $\varphi' \rightarrow \varphi$ and $\varphi \not\rightarrow \varphi'$, i.e., φ' holds on a proper subset of the set of words on which φ hold. In this case, our algorithm generates a word that satisfies φ and not φ' and adds it to N to eliminate φ .
- Third, $\varphi' \not\rightarrow \varphi$, i.e., φ' does not hold on a subset of the set of words on which φ hold. In this case, our algorithm generates a word w that satisfies φ' and not φ and adds it to N to eliminate φ' .

By handling the cases mentioned above, we obtain an algorithm with guarantees exactly same as the semi-symbolic algorithm in Section 4.1.

5 Experiments

In this section, we evaluate the performance of the proposed algorithms using three case studies. First, we evaluate the performance of Algorithm 1, referred to as SYM_{DFA} , and compare it to a baseline counterexample-guided algorithm by Avellaneda and Petrenko (2018), referred to as CEG_{DFA} (case study 1). Then, we evaluate the performance of the proposed semi-symbolic algorithm (Section 4.1), referred to as $\text{S-SYM}_{\text{LTL}}$, and counterexample-guided algorithm (Section 4.2), referred to as CEG_{LTL} , for learning LTL formulas (case studies 2 and 3).

In $\text{S-SYM}_{\text{LTL}}$, we fix the time horizon K to a natural number, instead of the double exponential theoretical upper bound of $2^{2^{n+1}}$. Using this heuristic means that $\text{S-SYM}_{\text{LTL}}$ does not solve Problem 2, but we demonstrate that we produce good enough formulas in practice.

In addition, we implement two existing heuristics from Avellaneda and Petrenko (2018) to all the algorithms. First, we implement in every algorithm, we learn models in an incremental manner, i.e., we start by learning DFAs (resp. LTL formulas) of size 1 and then the size by 1. We repeat the process till bound n . Second, we use a set of positive words P' instead of P that starts as an empty set, and at each iteration of the algorithm, if the inferred language does not contain some words from the original set P , we then extend P' with one of such words, preferably the shortest one. This last heuristic helps when dealing with large input samples, by using as fewest words from P .

We implemented every algorithm using python, using PySAT (Ignatiev, Morgado, and Marques-Silva, 2018) for learning DFA, and using an ASP (Baral, 2003) encoding that we solve using clingo (Gebser et al., 2017) for learning LTL formulas. Overall, we ran the experiments using 8 GiB of RAM and two CPU cores.

Learning DFAs For this case study, we consider a set of 28 random DFAs of size 2 to 10 generated using AALpy (Muškardin et al., 2022). Using each random DFA, we generate a set of 1000 positive words of length 1 to 10. We run algorithms CEG_{DFA} and SYM_{DFA} with a timeout $TO = 1000\text{s}$, and for n up to 10.

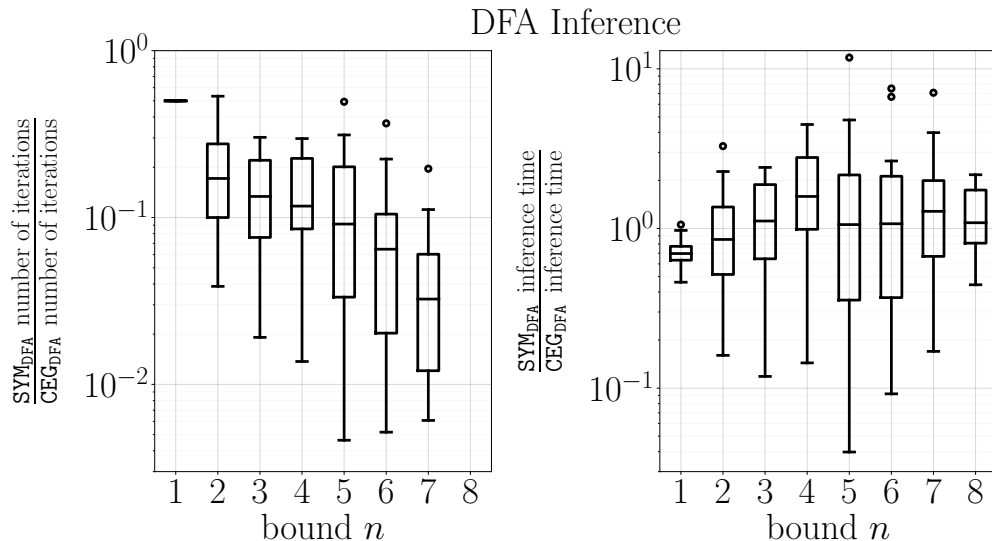


Fig. 2: Comparison of SYM_{DFA} and CEG_{DFA} in terms of the runtime and the number of iterations of the main loop.

Figure 2 shows a comparison between the performance of SYM_{DFA} and CEG_{DFA} in terms of the inference time and the required number of iterations of the main loop. On the left plot, the average ratio of the number of iterations is 0.14 which, in fact, shows that SYM_{DFA} required noticeably less number of iterations in comparison with CEG_{DFA} . On the right plot, the average ratio of the inference time is 1.09 which shows that the inference of the two algorithms is comparable and yet SYM_{DFA} is computationally less expensive due to requiring less iterations.

Learning Common LTL Patterns In this case study, we generate sample words using 12 common LTL patterns (Dwyer, Avrunin, and Corbett, 1999). For instance, $\psi_1 = \mathbf{G} a_0$, $\psi_2 = \mathbf{G}(a_1 \rightarrow \mathbf{G}(\neg a_0))$ and $\psi_3 = (\mathbf{G}(\neg a_0)) \vee (\mathbf{F}(a_0 \wedge (\mathbf{F}(a_1))))$. We refer the readers to appendix for a complete list of LTL patterns. Using each of these 12 ground truth LTL formulas, we generate a sample of 10000 positive words of length 10. Then, we infer LTL formulas for each sample using CEG_{LTL} and $\text{S-SYM}_{\text{LTL}}$, separately. For both algorithms, we set the maximum formula size $n = 10$ and a timeout of $TO = 1000\text{s}$, and for $\text{S-SYM}_{\text{LTL}}$, we set the time horizon $K = 8$. Figure 3 represents a comparison between the mentioned algorithms in terms of inference time for the ground truth LTL formulas ψ_1 , ψ_2 , and ψ_3 .

On average, $\text{S-SYM}_{\text{LTL}}$ is 173.9% faster than CEG_{LTL} for all the 12 samples. Our results show that LTL formulas inferred by $\text{S-SYM}_{\text{LTL}}$ are more or equally specific (i.e., included in term of language) than the ground truth LTL formulas for five samples out of the 12 samples, while the LTL formulas inferred by CEG_{LTL} are equally or more specific than the ground truth LTL formulas for three samples out of the 12 samples.

Learning LTL from Trajectories of Unmanned Aerial Vehicle (UAV) In this case study, we implement $\text{S-SYM}_{\text{LTL}}$ and CEG_{LTL} using sample words of a simulated unmanned aerial vehicle (UAV) for learning LTL formulas. Here, we use 10000 words clustered into three bundles using k -means clustering approach. Each words summarizes selective binary features such as $x_0 = 0$: “low battery” ($x_0 = 1$: “high battery”), $x_1 = 0$: “glide” ($x_1 = 1$: “thrust”), $x_2 = 0$: “change yaw angle” ($x_2 = 1$: “not change yaw angle”), $x_3 = 0$: “change roll angle” ($x_3 = 1$: “not change roll angle”), etc. We set $n = 10$, $K = 8$, and a timeout of $TO = 3600\text{s}$. We inferred LTL formulas for each cluster using CEG_{LTL} and $\text{S-SYM}_{\text{LTL}}$. Figure 4 depicts a comparison between CEG_{LTL} and $\text{S-SYM}_{\text{LTL}}$ in terms of the inference time for three clusters. Our results show that, on average, $\text{S-SYM}_{\text{LTL}}$ is 260.73% faster than CEG_{LTL} . Two examples of the learned LTL formulas from the UAV words are $\mathbf{F}(x_1) \vee \neg \mathbf{G}(x_1)$ which reads as “either the UAV always glides, or it never glides” and $\mathbf{G}(x_2 \rightarrow x_3)$ which reads as “a change in yaw angle is always accompanied by a change in roll angle”.

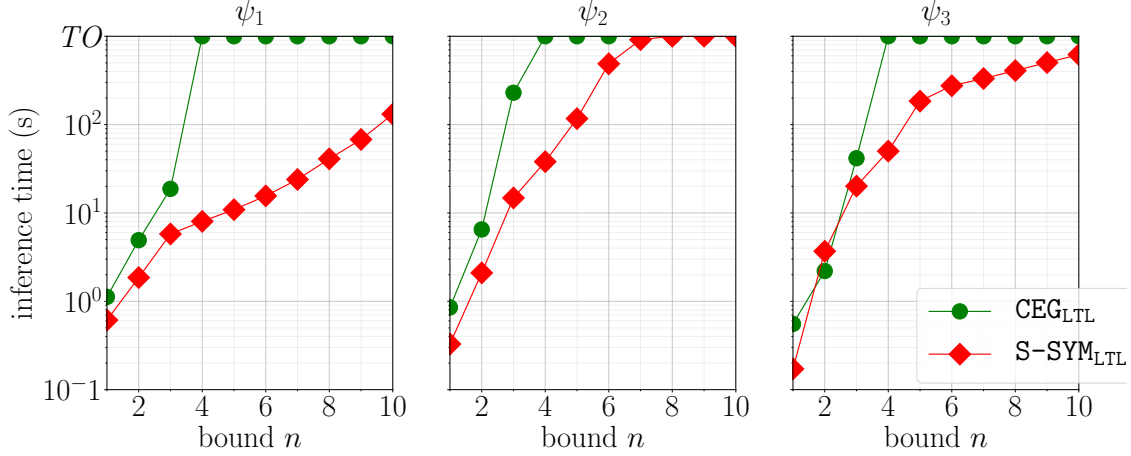


Fig. 3: Comparison of $S\text{-SYM}_{LTL}$ and CEG_{LTL} in terms of the inference time for three LTL ground truth formulas.

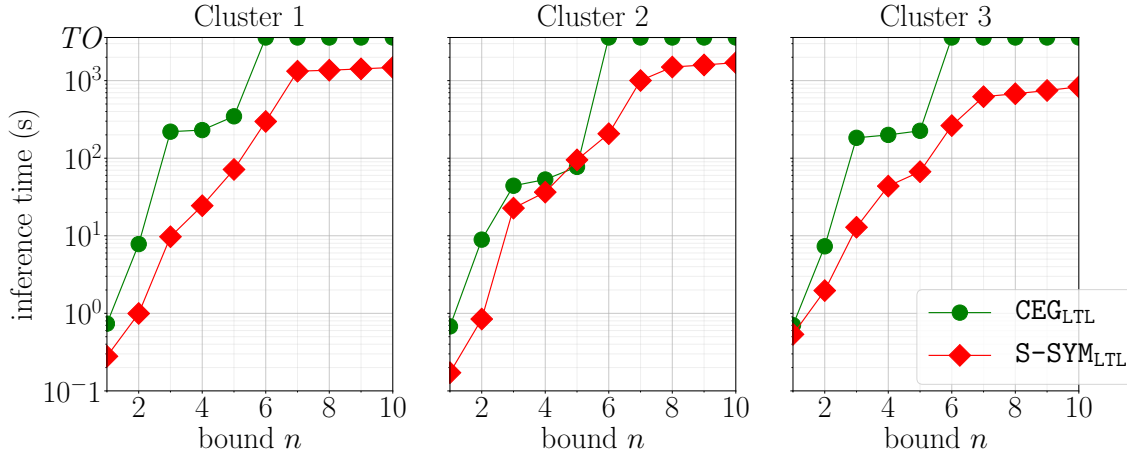


Fig. 4: Comparison of SYM_{DFA} and CEG_{DFA} in terms of the runtime for three clusters of words taken from a UAV.

6 Conclusion

We present novel algorithms for learning DFAs and LTL formulas from positive examples only. Our algorithms rely on conciseness and language minimality as regularizers to learn meaningful models. We demonstrate the efficacy of our algorithms on three case studies.

A natural direction of future work is to lift our techniques to tackle learning from positive examples for other finite state machines, e.g., non-deterministic finite automata (Sipser, 1997), and more expressive temporal logics, e.g., linear dynamic logic (LDL) (Giacomo and Vardi, 2013).

Acknowledgments

We are especially grateful to Dhananjay Raju for introducing us to Answer Set Programming and guiding us in using it to solve our SAT problem. This work has been supported by the Defense Advanced Research Projects Agency (DARPA) (Contract number HR001120C0032), Army Research Laboratory (ARL) (Contract number W911NF2020132)

and ACC-APG-RTP W911NF), National Science Foundation (NSF) (Contract number 1646522), and Deutsche Forschungsgemeinschaft (DFG) (Grant number 434592664).

References

- Angluin, D. 1980. Finding patterns common to a set of strings. *Journal of Computer and System Sciences*, 21(1): 46–62.
- Avellaneda, F.; and Petrenko, A. 2018. Inferring DFA without Negative Examples. In Unold, O.; Dyrka, W.; and Wieczorek, W., eds., *Proceedings of the 14th International Conference on Grammatical Inference, ICGI 2018, Wrocław, Poland, September 5-7, 2018*, volume 93 of *Proceedings of Machine Learning Research*, 17–29. PMLR.
- Baral, C. 2003. *Knowledge Representation, Reasoning and Declarative Problem Solving*. Cambridge University Press.
- Biermann, A. W.; and Feldman, J. A. 1972. On the Synthesis of Finite-State Machines from Samples of Their Behavior. *IEEE Trans. Computers*, 21(6): 592–597.
- Camacho, A.; and McIlraith, S. A. 2019. Learning Interpretable Models Expressed in Linear Temporal Logic. In *ICAPS*, 621–630. AAAI Press.
- Carrasco, R. C.; and Oncina, J. 1999. Learning deterministic regular grammars from stochastic samples in polynomial time. *RAIRO Theor. Informatics Appl.*, 33(1): 1–20.
- Chou, G.; Ozay, N.; and Berenson, D. 2022. Learning temporal logic formulas from suboptimal demonstrations: theory and experiments. *Auton. Robots*, 46(1): 149–174.
- Dwyer, M. B.; Avrunin, G. S.; and Corbett, J. C. 1999. Patterns in Property Specifications for Finite-State Verification. In Boehm, B. W.; Garlan, D.; and Kramer, J., eds., *Proceedings of the 1999 International Conference on Software Engineering, ICSE '99, Los Angeles, CA, USA, May 16-22, 1999*, 411–420. ACM.
- Ehlers, R.; Gavran, I.; and Neider, D. 2020. Learning Properties in $LTL \sqcap ACTL$ from Positive Examples Only. In *FMCAD*, 104–112. IEEE.
- Fandinno, J.; Laferrière, F.; Romero, J.; Schaub, T.; and Son, T. C. 2021. Planning with Incomplete Information in Quantified Answer Set Programming.
- Gebser, M.; Kaminski, R.; Kaufmann, B.; and Schaub, T. 2017. Multi-shot ASP solving with clingo. *CoRR*, abs/1705.09811.
- Giacomo, G. D.; and Vardi, M. Y. 2013. Linear Temporal Logic and Linear Dynamic Logic on Finite Traces. In *IJCAI*, 854–860. IJCAI/AAAI.
- Giacomo, G. D.; and Vardi, M. Y. 2015. Synthesis for LTL and LDL on Finite Traces. In *IJCAI*, 1558–1564. AAAI Press.
- Gold, E. M. 1967. Language Identification in the Limit. *Inf. Control.*, 10(5): 447–474.
- Gold, E. M. 1978. Complexity of Automaton Identification from Given Data. *Inf. Control.*, 37(3): 302–320.
- Grinchtein, O.; Leucker, M.; and Piterman, N. 2006. Inferring Network Invariants Automatically. In *IJCAR*, volume 4130 of *Lecture Notes in Computer Science*, 483–497. Springer.
- Gunning, D.; Stefik, M.; Choi, J.; Miller, T.; Stumpf, S.; and Yang, G.-Z. 2019. XAI2014; Explainable artificial intelligence. *Science Robotics*, 4(37): eaay7120.
- Heule, M.; and Verwer, S. 2010. Exact DFA Identification Using SAT Solvers. In *10th International Colloquium of Grammatical Inference: Theoretical Results and Applications, ICGI '10*, volume 6339 of *LNCS*, 66–79. Springer.
- Ignatiev, A.; Morgado, A.; and Marques-Silva, J. 2018. PySAT: A Python Toolkit for Prototyping with SAT Oracles. In *SAT*, 428–437.
- Jha, S.; Tiwari, A.; Seshia, S. A.; Sahai, T.; and Shankar, N. 2019. TeLEx: learning signal temporal logic from positive examples using tightness metric. *Formal Methods Syst. Des.*, 54(3): 364–387.
- Kasenberg, D.; and Scheutz, M. 2017. Interpretable apprenticeship learning with temporal logic specifications. In *CDC*, 4914–4921. IEEE.
- Lemieux, C.; Park, D.; and Beschastnikh, I. 2015. General LTL Specification Mining (T). In *ASE*, 81–92. IEEE Computer Society.
- Li, C.; and Manyà, F., eds. 2021. *Theory and Applications of Satisfiability Testing - SAT 2021 - 24th International Conference, Barcelona, Spain, July 5-9, 2021, Proceedings*, volume 12831 of *Lecture Notes in Computer Science*. Springer.
- Molnar, C. 2022. *Interpretable Machine Learning*. 2 edition.

- Muškardin, E.; Aichernig, B.; Pill, I.; Pferscher, A.; and Tappler, M. 2022. AALpy: an active automata learning library. *Innovations in Systems and Software Engineering*, 1–10.
- Neider, D.; and Gavran, I. 2018. Learning Linear Temporal Properties. In Bjørner, N.; and Gurfinkel, A., eds., *2018 Formal Methods in Computer Aided Design, FMCAD 2018, Austin, TX, USA, October 30 - November 2, 2018*, 1–10. IEEE.
- Pnueli, A. 1977. The Temporal Logic of Programs. In *18th Annual Symposium of Foundations of Computer Science, FOCS '77*, 46–57. IEEE Computer Society.
- Rabin, M.; and Scott, D. 1959. Finite Automata and Their Decision Problems. *IBM Journal of Research and Development*, 3: 114–125.
- Raha, R.; Roy, R.; Fijalkow, N.; and Neider, D. 2022. Scalable Anytime Algorithms for Learning Fragments of Linear Temporal Logic. In *TACAS (1)*, volume 13243 of *Lecture Notes in Computer Science*, 263–280. Springer.
- Roy, R.; Fisman, D.; and Neider, D. 2020. Learning Interpretable Models in the Property Specification Language. In *IJCAI*, 2213–2219. ijcai.org.
- Royal-Society. 2019. Explainable AI: The Basics.
- Shvo, M.; Li, A. C.; Icarte, R. T.; and McIlraith, S. A. 2021. Interpretable Sequence Classification via Discrete Optimization. In *AAAI*, 9647–9656. AAAI Press.
- Sipser, M. 1997. *Introduction to the theory of computation*. PWS Publishing Company.
- Sistla, A. P.; and Clarke, E. M. 1985. The Complexity of Propositional Linear Temporal Logics. *J. ACM*, 32(3): 733–749.
- Stolcke, A.; and Omohundro, S. 1992. Hidden Markov Model Induction by Bayesian Model Merging. In Hanson, S.; Cowan, J.; and Giles, C., eds., *Advances in Neural Information Processing Systems*, volume 5. Morgan-Kaufmann.
- Vazquez-Chanlatte, M.; Jha, S.; Tiwari, A.; Ho, M. K.; and Seshia, S. A. 2018. Learning Task Specifications from Demonstrations. In *NeurIPS*, 5372–5382.
- Weiss, G.; Goldberg, Y.; and Yahav, E. 2018. Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, 5244–5253. PMLR.
- Zhu, S.; Tabajara, L. M.; Li, J.; Pu, G.; and Vardi, M. Y. 2017. Symbolic LTLf Synthesis. In *IJCAI*, 1362–1369. ijcai.org.

A Results for the Symbolic Algorithm in Section 3

A.1 The Theoretical Guarantees

Before we prove the theoretical guarantees of the algorithm (Theorem 1), we re-state the constraints used for the symbolic algorithm (Algorithm 1) for convenience. We begin with the constraints for Φ_{DFA} :

$$\bigwedge_{p \in [n]} \bigwedge_{a \in \Sigma} \left[\bigvee_{q \in [n]} d_{p,a,q} \wedge \bigwedge_{q \neq q' \in [n]} [\neg d_{p,a,q} \vee \neg d_{p,a,q'}] \right] \quad (3)$$

We then provide the constraints for Φ_P as follows:

$$x_{\varepsilon,1} \wedge \bigwedge_{q \in \{2, \dots, n\}} \neg x_{\varepsilon,q} \quad (4)$$

$$\bigwedge_{ua \in \text{Pref}(P)} \bigwedge_{p,q \in [n]} [x_{u,p} \wedge d_{p,a,q}] \rightarrow x_{ua,q} \quad (5)$$

$$\bigwedge_{w \in P} \bigwedge_{q \in [n]} x_{w,q} \rightarrow f_q \quad (6)$$

Next, we provide the constraints for $\Phi_{\subseteq \mathcal{A}}$

$$y_{1,1}^{\mathcal{A}} \quad (7)$$

$$\bigwedge_{q=\delta(p,a)} \bigwedge_{p',q' \in [n]} \bigwedge_{a \in \Sigma} \left[[y_{p,p'}^{\mathcal{A}} \wedge d_{p',a,q'}] \rightarrow y_{q,q'}^{\mathcal{A}} \right] \quad (8)$$

$$\bigwedge_{p \notin F} \bigwedge_{p' \in [n]} \left[y_{p,p'}^{\mathcal{A}} \rightarrow \neg f_{p'} \right] \quad (9)$$

Finally, we end with the constraints for $\Phi_{\supseteq \mathcal{A}}$:

$$z_{0,1,1} \quad (10)$$

$$\bigwedge_{i \in [n^2]} \left[\bigvee_{q,q' \in [n]} z_{i,q,q'} \wedge \left[\bigwedge_{\substack{p \neq q \in [n] \\ p' \neq q' \in [n]}} \neg z_{i,p,p'} \vee \neg z_{i,q,q'} \right] \right] \quad (11)$$

$$\bigwedge_{\substack{p,q \in [n] \\ p',q' \in [n]}} \left[[z_{i,p,p'} \wedge z_{i+1,q,q'}] \rightarrow \bigvee_{\substack{a \in \Sigma \text{ where} \\ q=\delta(p,a)}} d_{p',a,q'} \right] \quad (12)$$

$$\bigvee_{i \in [n^2]} \bigvee_{\substack{q \in [n] \\ q' \in F}} \left[z_{i,q,q'} \wedge \neg f_{q'} \right] \quad (13)$$

We now prove the correctness of the propositional formula $\Phi^{\mathcal{A}}$ that we construct using the above constraints:

Theorem 3. *Let $\Phi^{\mathcal{A}}$ be as defined above. Then, we have the following:*

- (1) *If $\Phi^{\mathcal{A}}$ is satisfiable, then there exists a DFA \mathcal{A}' that is an n -description and $L(\mathcal{A}') \subset L(\mathcal{A})$.*
- (2) *If there exists a DFA \mathcal{A}' that is an n -description and $L(\mathcal{A}') \subset L(\mathcal{A})$, then $\Phi^{\mathcal{A}}$ is satisfiable.*

To prove the above theorem, we propose intermediate claims, all of which we prove first. In what follows, we assume v to be a model of $\Phi^{\mathcal{A}}$, $\mathcal{A}' = (Q', \Sigma, \delta', q_I', F')$ to be the DFA constructed from a model v of $\Phi^{\mathcal{A}}$ and $\mathcal{A} = (Q, \Sigma, \delta, q_I, F)$ to be the current hypothesis.

Claim. For all $u \in \text{Pref}(P)$, $\mathcal{A}' : q_I' \xrightarrow{u} q$ implies $v(x_{u,q}) = 1$.

Proof. We prove the claim using an induction on the length $|u|$ of the word u .

Base case: Let $u = \epsilon$. Based on the definition of runs, $\mathcal{A}' : q_I' \xrightarrow{\epsilon} q$ implies $q = q_I'$. Moreover, using Constraint 4, we have $v(x_{\epsilon,q}) = 1$ if and only if $q = q_I'$. These two facts prove the claim for the base case.

Induction hypothesis: As induction hypothesis, let $\mathcal{A}' : q_I' \xrightarrow{u} q$ implies $v(x_{u,q}) = 1$ for all words $u \in \text{Pref}(P)$ of length $\leq k$. Now, consider the run $\mathcal{A}' : q_I' \xrightarrow{u} q \xrightarrow{a} q'$ for some $ua \in \text{Pref}(P)$. For this run, based on the induction hypothesis and the construction of \mathcal{A}' , $v(x_{u,q}) = 1$ and $v(d_{p,a,q}) = 1$. Now, using Constraint 5, $v(x_{u,q}) = 1$ and $v(d_{p,a,q}) = 1$ implies $v(x_{ua,q}) = 1$, proving the claim.

Claim. For all $w \in \Sigma^*$, $\mathcal{A} : q_I \xrightarrow{w} q$ and $\mathcal{A}' : q_I' \xrightarrow{w} q'$ imply $v(y_{q,q'}^{\mathcal{A}}) = 1$.

Proof. We prove this using induction of the length $|w|$ of the word w .

Base case: Let $w = \epsilon$. Based on the definition of runs, $\mathcal{A} : q_I \xrightarrow{\epsilon} q$, $\mathcal{A}' : q_I' \xrightarrow{\epsilon} q'$ implies $q = q_I$ and $q = q_I'$. Moreover, using Constraint 7, $q = q_I$ and $q' = q_I'$ implies $v(y_{q,q'}^{\mathcal{A}}) = 1$. These two facts prove the claim for the base case.

Induction hypothesis: As induction hypothesis, let $\mathcal{A} : q_I \xrightarrow{w} q$ and $\mathcal{A}' : q_I' \xrightarrow{w} q'$ imply $v(y_{q,q'}^{\mathcal{A}}) = 1$ for all words $w \in \Sigma^*$ of length $\leq k$. Now, consider the runs $\mathcal{A} : q_I \xrightarrow{w} p \xrightarrow{a} q$ and $\mathcal{A}' : q_I' \xrightarrow{w} p' \xrightarrow{a} q'$ for some word $wa \in \Sigma^*$. For this run, based on the induction hypothesis and the construction of \mathcal{A}' , $v(y_{p,p'}^{\mathcal{A}}) = 1$ and $v(d_{p',a,q'}) = 1$. Now, using Constraint 8, together $v(y_{p,p'}^{\mathcal{A}}) = 1$ and $v(d_{p',a,q'}) = 1$ imply $v(y_{q,q'}^{\mathcal{A}}) = 1$ (where $q = \delta(p, a)$), proving the claim.

Claim. $v(z_{i,q,q'}) = 1$ implies there exists $w \in \Sigma^i$ with runs $\mathcal{A} : q_I \xrightarrow{w} q$ and $\mathcal{A}' : q_{I'} \xrightarrow{w} q'$.

Proof. We proof this using an induction on the parameter i .

Base case: Let $i = 0$. Based on the Constraints 10 and 11, $v(z_{0,q,q'}) = 1$ implies $q = q_I$ and $q = q_{I'}$. Now, there always exists a word of length 0, i.e., $w = \epsilon$, for which $\mathcal{A} : q_I \xrightarrow{\epsilon} q$ and $\mathcal{A}' : q_{I'} \xrightarrow{\epsilon} q'$ proving the claim for the base case.

Induction hypothesis: As induction hypothesis, let $v(z_{k,p,p'}) = 1$ hold and thus, w be a word of length k such that $\mathcal{A} : q_I \xrightarrow{w} p$ and $\mathcal{A}' : q_{I'} \xrightarrow{w} p'$. Now, consider $v(z_{k+1,q,q'}) = 1$. Based on Constraint 12, for some $a \in \Sigma$ such that $q = \delta(p, a)$, $v(d_{p',a,q'}) = 1$. Thus, on the word wa , we can have the runs $\mathcal{A} : q_I \xrightarrow{w} p \xrightarrow{a} q$ and $\mathcal{A}' : q_{I'} \xrightarrow{w} p' \xrightarrow{a} q'$, proving the claim.

We are now ready to prove Theorem 3, i.e., the correctness of $\Phi^{\mathcal{A}}$.

Proof (of Theorem 3). For the first direction, let consider $\Phi^{\mathcal{A}}$ is satisfiable and v is a model of $\Phi^{\mathcal{A}}$ and \mathcal{A}' be the DFA constructed from the model v . First, based on Claim A.1, for all $w \in P$, $\mathcal{A}' : q_{I'} \xrightarrow{w} q$. This further implies $q \in F'$, using Constraint 6. Thus, \mathcal{A}' accepts all positive words and hence, is an n -description. Second, based on Claim A.1, for all words $w \in \Sigma^*$, $\mathcal{A} : q_I \xrightarrow{w} q$ and $\mathcal{A}' : q_{I'} \xrightarrow{w} q'$ imply $v(y_{q,q'}^{\mathcal{A}}) = 1$. Thus, if $q \notin F$, then $q' \notin F$, using Constraint 9, implying $L(\mathcal{A}') \subseteq L(\mathcal{A})$. Third, based on Claim A.1, $v(z_{i,q,q'}) = 1$ implies there exists $w \in \Sigma^i$ with runs $\mathcal{A} : q_I \xrightarrow{w} q$ and $\mathcal{A}' : q_{I'} \xrightarrow{w} q'$. Combining this with Constraint 13, we deduce that there exists $w \in \Sigma^*$ with length $\leq n^2$ with a run of \mathcal{A} ending in F and run of \mathcal{A}' not ending in F' . This shows that $L(\mathcal{A}) \neq L(\mathcal{A}')$. We thus conclude that \mathcal{A} is an n -description and $L(\mathcal{A}') \subset L(\mathcal{A})$.

For the other direction, based on a suitable DFA \mathcal{A}' , we construct an assignment v for all the introduce propositional variables. First, we set $v(d_{p,a,q}) = 1$ if $\delta'(p, a) = q$ and $v(f_q) = 1$ if $q \in F'$. Since δ is a deterministic function, it is clear that v satisfies the Constraint 3. Similarly, we set $v(x_{u,q}) = 1$ if $\mathcal{A}' : q_{I'} \xrightarrow{u} q$ for some $u \in Pref(P)$. It is again a simple exercise to check that v satisfies Constraints 4 to 6. Next, we set $v(y_{q,q'}^{\mathcal{A}}) = 1$ if there are runs $\mathcal{A} : q_I \xrightarrow{w} q$ and $\mathcal{A}' : q_{I'} \xrightarrow{w} q'$ on some word $w \in \Sigma^*$. The runs of the DFAs can be simulated using simple breadth-first search techniques. Again, with a quick check, one can see that model v satisfies Constraints 7 to 9. Finally, we set assignment to $z_{i,q,q'}$ exploiting a word w which permits runs $\mathcal{A} : q_I \xrightarrow{w} q$ and $\mathcal{A}' : q_{I'} \xrightarrow{w} q'$, where q is in F but q' not in F' . In particular, we set $v(z_{i,q,q'}) = 1$ for $i = |u|$ and $\mathcal{A} : q_I \xrightarrow{u} q$ and $\mathcal{A}' : q_{I'} \xrightarrow{u} q'$ for all prefixes u of w . Clearly, such an assignment encodes a synchronized run of the DFAs \mathcal{A} and \mathcal{A}' on word w and thus, satisfies Constraints 10 to 13.

The correctness of the encoding $\Phi^{\mathcal{A}}$ provides in a straightforward manner the correctness of the symbolic algorithm, since the search space of DFA of upto size n is finite.

A.2 The complete algorithm with heuristics

We first present the complete symbolic algorithm (i.e., along with the main heuristics) for solving the OCC problem for DFA (Problem 1). The pseudocode is sketched in Algorithm 2. Compared to the algorithm presented in Algorithm 1, we make a few modifications to improve performance. First, we introduce a set P' (also described in Section 5) to store the set of positive words necessary for learning the hypothesis DFA \mathcal{A}' . Second, we incorporate an incremental DFA learning, meaning that we search for DFAs satisfying the propositional formula $\Phi^{\mathcal{A}}$ of increasing size (starting from size 1). To reflect this, we extend $\Phi^{\mathcal{A}}$ with the size parameter, represented using $\Phi^{\mathcal{A},m}$, to search for a DFA of size m .

A.3 Comparison with a Semi-symbolic Approach

We have introduced counter-example guided, semi-symbolic and symbolic approaches in this paper. Our exploration of these methods wouldn't be complete if we didn't also tried a semi-symbolic algorithm for learning DFA. Hence, we introduce S-SYM_{DFA}, a semi-symbolic approach for learning DFA. This is done in a similar fashion than for LTL (Algorithm 3), but with DFA instead. Hence, we define the SAT problem as $\Phi^{\mathcal{A}} := \Phi_{\text{DFA}} \wedge \Phi_P \wedge \Phi_N \wedge \Phi_{\neg \mathcal{A}}$. In practice, S-SYM_{DFA} is always worse than SYM_{DFA}, both in term of inference time (in average, 3.2 times more) and number of iterations (in average, 2.1 times more), as demonstrated in Figure 5.

Algorithm 2: Symbolic Algorithm for Learning DFA (with heuristics)

Input: Positive words P , bound n

- 1: $\mathcal{A} \leftarrow \mathcal{A}_{\Sigma^*}, \Phi^{\mathcal{A}} \leftarrow \Phi_{\text{DFA}} \wedge \Phi_P$
- 2: $P' \leftarrow \emptyset, m \leftarrow 1$
- 3: **while** $m \leq n$ **do**
- 4: $\Phi^{\mathcal{A},m} \leftarrow \Phi_{\text{DFA}}^m \wedge \Phi_{P'} \wedge \Phi_{\subseteq \mathcal{A}} \wedge \Phi_{\not\subseteq \mathcal{A}}$
- 5: **if** no model v satisfies $\Phi^{\mathcal{A},m}$ **then**
- 6: $m \leftarrow m + 1$
- 7: **else**
- 8: $\mathcal{A}' \leftarrow$ DFA constructed from v
- 9: **if** exists $w \in P \setminus L(\mathcal{A}')$ **then**
- 10: Add the shortest of such w to P'
- 11: **else**
- 12: $\mathcal{A} \leftarrow \mathcal{A}'$
- 13: **end if**
- 14: **end if**
- 15: **end while**
- 16: **return** \mathcal{A}

Algorithm 3: Semi-symbolic Algorithm for LTL formulas

Input: Positive words P , bound n

- 1: $N \leftarrow \emptyset, D \leftarrow \emptyset$
- 2: $\varphi \leftarrow \varphi_{\Sigma^*}, \Omega^{N,D} := \Psi_{\text{LTL}} \wedge \Psi_P$
- 3: **while** $\Omega^{N,D}$ is satisfiable (with model v) **do**
- 4: $\varphi' \leftarrow$ LTL formula constructed from v
- 5: **if** $\varphi' \rightarrow \varphi$ **then**
- 6: Update φ to φ'
- 7: **else**
- 8: Add w to N , where $w \models \neg\varphi' \wedge \varphi$
- 9: **end if**
- 10: $\Omega^{N,D} := \Psi_{\text{LTL}} \wedge \Psi_P \wedge \Psi_N \wedge \Psi_{\neq \varphi}$
- 11: **end while**
- 12: **return** φ

B Results for the Semi-symbolic Algorithm in Section 4.1

B.1 The Algorithm

We first present the pseudo-code of the semi-symbolic algorithm for solving the OCC problem for LTL formulas (Problem 2).

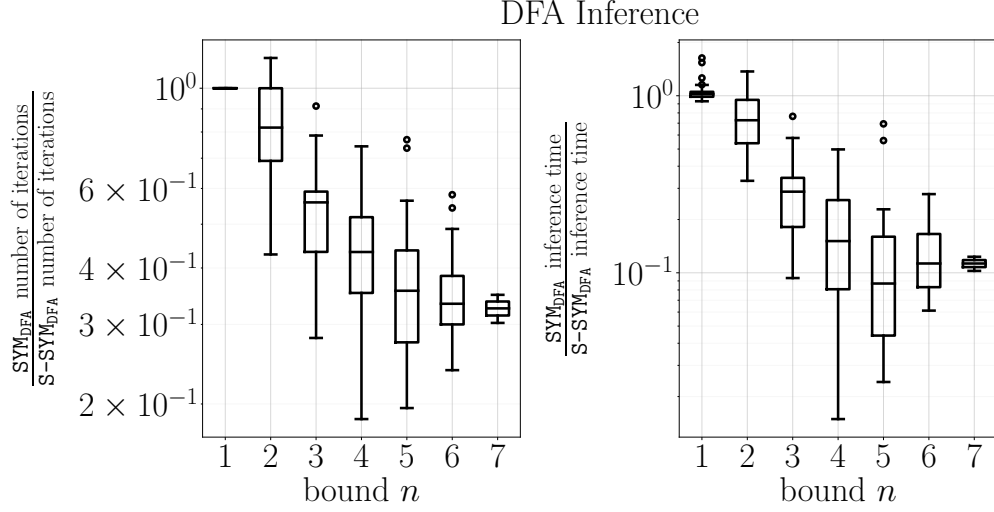


Fig. 5: Comparison of SYM_{DFA} and $\text{S-SYM}_{\text{DFA}}$ in terms of the inference time and the number of iterations of the main loop.

B.2 All the constraints

We here provide all the constraints necessary for constructing the main propositional formula $\Psi^{\varphi, N}$. We begin by listing the structural constraints required for Ψ_{LTL} .

$$\bigwedge_{i \in [n]} \left[\bigvee_{\lambda \in A} x_{i, \lambda} \wedge \bigwedge_{\lambda \neq \lambda' \in A} [\neg x_{i, \lambda} \vee \neg x_{i, \lambda'}] \right] \quad (14)$$

$$\bigwedge_{i \in \{2, \dots, n\}} \left[\bigvee_{j \in [i-1]} l_{i, j} \wedge \bigwedge_{j \neq j' \in [i-1]} [\neg r_{i, j} \vee \neg l_{i, j'}] \right] \quad (15)$$

$$\bigwedge_{i \in \{2, \dots, n\}} \left[\bigvee_{j \in [i-1]} r_{i, j} \wedge \bigwedge_{j \neq j' \in [i-1]} [\neg r_{i, j} \vee \neg r_{i, j'}] \right] \quad (16)$$

$$\bigvee_{a \in \Sigma} x_{1, a} \quad (17)$$

We then move on to the semantic constraints required for Ψ_P and Ψ_N .

$$\bigwedge_{i \in [n]} \bigwedge_{a \in \Sigma} x_{i,a} \rightarrow \left[\bigwedge_{t \in [|w|]} \begin{cases} y_{w,t}^i & \text{if } a = w[t] \\ \neg y_{w,t}^i & \text{if } a \neq w[t] \end{cases} \right] \quad (18)$$

$$\bigwedge_{\substack{i \in [n] \\ j, j' \in [i-1]}} x_{i,\vee} \wedge l_{i,j} \wedge r_{i,j'} \rightarrow \left[\bigwedge_{\tau \in [|w|]} \left[y_{w,t}^i \leftrightarrow y_{w,t}^j \vee y_{w,t}^{j'} \right] \right] \quad (19)$$

$$\bigwedge_{\substack{i \in [n] \\ j \in [i-1]}} x_{i,\mathbf{X}} \wedge l_{i,j} \rightarrow \left[\bigwedge_{t \in [|w|-1]} y_{w,t}^i \leftrightarrow y_{w,t+1}^j \right] \wedge \neg y_{w,|w|}^j \quad (20)$$

$$\bigwedge_{\substack{i \in [n] \\ j, j' \in [i]}} x_{i,\mathbf{U}} \wedge l_{i,j} \wedge r_{i,j'} \rightarrow \left[\bigwedge_{t \in [|w|]} \left[y_{w,t}^i \leftrightarrow \bigvee_{t \leq t' \in [|w|]} \left[y_{w,t'}^{j'} \wedge \bigwedge_{t \leq \tau < t'} y_{w,\tau}^j \right] \right] \right] \quad (21)$$

As mentioned in the Section 4.1, Ψ_P additionally imposes $\bigvee_{w \in P} y_{w,0}^n$ and Ψ_N imposes $\bigvee_{w \in N} \neg y_{w,0}^n$. Next, we provide the constraints in Ψ_{word} for symbolically encoding a word of length at most K .

$$\bigwedge_{t \in [K]} \left[\bigvee_{a \in \Sigma \cup \{\epsilon\}} p_{t,a} \wedge \bigwedge_{a \neq a' \in \Sigma \cup \{\epsilon\}} [\neg p_{t,a} \vee \neg p_{t,a'}] \right] \quad (22)$$

$$\bigwedge_{t \in [K-1]} p_{t,\epsilon} \rightarrow p_{t+1,\epsilon} \quad (23)$$

We now mention the constraints in Ψ_{sem} to ensure the meaning of the variables $z_{u,t}^{\varphi',i}$:

$$\bigwedge_{i \in [n]} \bigwedge_{a \in \Sigma} x_{i,a} \rightarrow \left[\bigwedge_{t \in [K]} z_{u,t}^{\varphi',i} \leftrightarrow p_{t,a} \right] \quad (24)$$

$$\bigwedge_{\substack{i \in [n] \\ j, j' \in [i-1]}} x_{i,\vee} \wedge l_{i,j} \wedge r_{i,j'} \rightarrow \left[\bigwedge_{t \in [K]} \left[z_{w,t}^{\varphi',i} \leftrightarrow z_{u,t}^{\varphi',j} \vee z_{u,t}^{\varphi',j'} \right] \right] \quad (25)$$

$$\bigwedge_{\substack{i \in [n] \\ j \in [i-1]}} x_{i,\mathbf{X}} \wedge l_{i,j} \rightarrow \left[\bigwedge_{t \in [K-1]} z_{w,t}^{\varphi',i} \leftrightarrow z_{u,t+1}^{\varphi',j} \right] \wedge \neg z_{u,|w|}^{\varphi',j} \quad (26)$$

$$\bigwedge_{\substack{i \in [n] \\ j, j' \in [i]}} x_{i,\mathbf{U}} \wedge l_{i,j} \wedge r_{i,j'} \rightarrow \left[\bigwedge_{t \in [K]} \left[z_{u,t}^{\varphi',i} \leftrightarrow \bigvee_{t \leq t' \in [K]} \left[z_{u,t'}^{\varphi',j'} \wedge \bigwedge_{t \leq \tau < t'} z_{u,\tau}^{\varphi',j} \right] \right] \right] \quad (27)$$

As evident, the above constraints are identical to the Constraints 18 to 21 imposed on variables $y_{w,t}^i$. The constraints on variables $z_{w,t}^{\varphi',i}$ are quite similar. The only difference is that, for hypothesis φ , we know the syntax DAG exactly and thus, can exploit it instead of using an encoding of the syntax DAG.

Finally as stated in Section 4.1, $\Psi_{\neq\varphi}$ is expressed as follows:

$$\Psi_{\text{word}} \wedge \Psi_{\text{sem}} \wedge z_{u,0}^{\varphi,n} \wedge \neg z_{u,0}^{\varphi',n} \quad (28)$$

We now prove the correctness of the encoding $\Psi^{\varphi,N}$ described using the constraints above.

Theorem 4. *Let $\Psi^{\varphi,N}$ be as defined above. Then, we have the following:*

- (1) *If $\Psi^{\varphi,N}$ is satisfiable, then there exists an LTL formula φ' that is an n -description, φ' does not hold on $w \in N$ and $\varphi \not\rightarrow \varphi'$.*
- (2) *If there exists a LTL formula φ' that is an n -description, φ' does not hold on $w \in N$ and $\varphi \not\rightarrow \varphi'$, then $\Psi^{\varphi,N}$ is satisfiable.*

To prove this theorem, we rely on intermediate claims, which we prove first. In what follows, v is a model of $\Psi^{\varphi,N}$, φ' is the LTL formula constructed from v and φ is the current hypothesis LTL formula.

Claim. For all $w \in P \cup N$, $v(y_{w,t}^i) = 1$ if and only if $w[t:] \models \varphi'[i]$.

The proof proceed via a structural induction over $\varphi'[i]$. For a proof of this claim, we refer the readers to the correctness proof of the encoding used by Neider and Gavran (2018).

Claim. $v(z_{w,t}^{\varphi',i}) = 1$ (resp. $v(z_{w,t}^{\varphi,i}) = 1$) if and only for a word w , $w[t:] \models \varphi'$ (resp. $w[t:] \models \varphi$).

The proof again proceeds via a structural induction on φ' (resp. φ) (similar to the one show by Neider and Gavran (2018)).

We are now ready to prove Theorem 4, i.e., the correctness of $\Omega^{N,D}$.

Proof (of Theorem 4). For the first direction, let consider $\Psi^{\varphi,N}$ is satisfiable and v is a model of $\Psi^{\varphi,N}$ and φ' be the LTL formula constructed from the model v . First, based on Claim B.2, we have that $v(y_{w,t}^n) = 1$ if and only if $w[t:] \models \varphi'$. Now, based on the constraints Ψ_P and Ψ_N , we observe that $v(y_{w,0}^n) = 1$ for all words $w \in P$ and $v(y_{w,0}^n) = 0$ for all words $w \in N$. Thus, combining the two observations, we conclude $w \models \varphi'$ for $w \in P$ and $w \not\models \varphi'$ for $w \in N$. Next, using Claim B.2 and conjunct $\Psi_{\neq\varphi}$, we conclude that there exists a word w on which φ holds and φ' does not. Thus, in total, we obtain φ' to be an n -description which does not hold on $w \in N$ and $\varphi \not\rightarrow \varphi'$.

For the other direction, based on a suitable LTL formula φ' , we construct an assignment v for all the introduce propositional variables. First, we set $v(x_{i,\lambda}) = 1$ if Node i is labeled with operator λ and $v(l_{i,j}) = 1$ (resp. $v(r_{i,j}) = 1$) if left (resp. right) child of Node i is Node j . Since φ is a valid LTL formula, it is clear that the structural constraints will be satisfied by v . Similarly, we set $v(y_{w,t}^i) = 1$ if $w[t:] \models \varphi'[i]$ for some $w \in P \cup N$, $t \in [|w|]$. It is again a simple exercise to check that v satisfies Constraints 18 to 21. Next, we set $v(z_{w,t}^{\varphi}) = 1$ and $v(z_{w,t}^{\varphi',n})$ for a word w for which $w \models \varphi$ and $w \not\models \varphi'$. It is again easy to check that v satisfies $\Psi_{\neq\varphi}$.

The correctness of the encoding $\Psi^{\varphi,N}$ provides in a straightforward manner the correctness of the semi-symbolic algorithm, since the search space of LTL formula of upto size n is finite.

C A symbolic algorithm for learning LTL formulas

We now describe few modifications to the semi-symbolic algorithm presented in Section 4.1 to convert it into a *completely* symbolic approach. This algorithm relies entirely on the hypothesis LTL formula φ for constructing a propositional formula Ψ^φ that guides the search of the next hypothesis. Precisely, the formula Ψ^φ has the properties that: (1) Ψ^φ is satisfiable if and only if there exists an LTL formula φ' that is an n -description and $\varphi \not\rightarrow \varphi'$ and $\varphi' \rightarrow \varphi$; and (2) based on a model v of $\Psi^{\varphi,N}$, one can construct such an LTL formula.

The algorithm, sketched in Algorithm 4, follows the same framework as Algorithm 1. We here make necessary modifications to search for an LTL formula. Also, the propositional formula Ψ^φ has a construction similar to $\Omega^{N,D}$, with the exception that Ψ_N is replaced by $\Psi_{\rightarrow\varphi}$.

We here only describe the construction of the conjunct $\Psi_{\rightarrow\varphi}$ which reuses the variables and constraints already introduced in Section 4.1.

$$\Psi_{\rightarrow\varphi} := \forall_{t \in [K], a \in \Sigma} p_{t,a} : \left[\Psi_{\text{word}} \wedge \Psi_{\text{sem}} \right] \rightarrow z_{w,t}^{\varphi',n} \rightarrow z_{w,t}^{\varphi,n} \quad (29)$$

Intuitively, the above constraint says that if for all words u of length $\leq K$, if φ' holds on u , then so must φ .

Algorithm 4: Symbolic Algorithm for Learning LTL

Input: Positive words P , bound n

- 1: $\varphi \leftarrow \varphi_{\Sigma^*}, \Psi^\varphi := \Psi_{\text{LTL}} \wedge \Psi_P$
 - 2: **while** Ψ^φ is satisfiable (with model v) **do**
 - 3: $\varphi \leftarrow$ LTL formula constructed from v
 - 4: $\Psi^\varphi := \Psi_{\text{LTL}} \wedge \Psi_P \wedge \Psi_{\rightarrow\varphi} \wedge \Psi_{\not\rightarrow\varphi}$
 - 5: **end while**
 - 6: **return** φ
-

C.1 Evaluation of SYM_{LTL}

We refer to this symbolic algorithm for learning LTL formulas (Algorithm 4) as SYM_{LTL} . We implement SYM_{LTL} using QASP2QBF Fandinno et al. (2021). SYM_{LTL} has an inference time several orders of magnitude above the inference time of $\text{S-SYM}_{\text{LTL}}$, as demonstrated in Figure 6. This can be explained by the choice of the solver, and the inherent complexity of the problem due to quantifiers. On the third experiment (Section 5), SYM_{LTL} timed out even for $n = 1$.

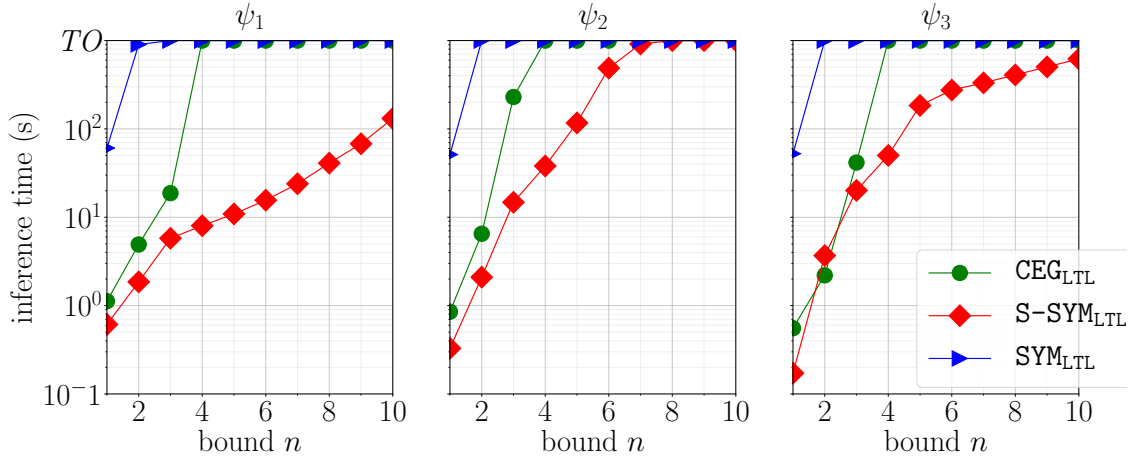


Fig. 6: For each sample of the second experiment (Section 5), comparison of the inference time between CEG_{LTL} , $\text{S-SYM}_{\text{LTL}}$ and SYM_{LTL} .

D The Counter-example guided Algorithm for learning LTL formula

We present the full pseudo-code for the counter-example guided algorithm from Section 4.2 in Algorithm 5. This algorithm is also referred as CEG_{LTL} .

E LTL patterns

We use ground-truth LTL formulas from Table 1 to generate the samples for the second experiment (Section 5). These LTL formulas are categorized in different categories, i.e., absence, existence, universality, and disjunctions of common LTL patterns.

Algorithm 5: CEG Algorithm for LTL formulas

Input: Positive words P , bound n

```

1:  $N \leftarrow \emptyset, D \leftarrow \emptyset$ 
2:  $\varphi \leftarrow \varphi_{\Sigma^*}, \Omega^{N,D} := \Psi_{\text{LTL}} \wedge \Psi_P$ 
3: while  $\Omega^{N,D}$  is satisfiable (with model  $v$ ) do
4:    $\varphi' \leftarrow \varphi^v$ 
5:   if  $\varphi' \leftrightarrow \varphi$  then
6:     Add  $\varphi'$  to  $D$ 
7:   else
8:     if  $\varphi' \rightarrow \varphi$  then
9:       Add  $w$  to  $N$ , where  $w \models \neg\varphi \wedge \varphi'$ 
10:       $\varphi \leftarrow \varphi'$ 
11:    else
12:      Add  $w$  to  $N$ , where  $w \models \neg\varphi' \wedge \varphi$ 
13:    end if
14:  end if
15:   $\Omega^{N,D} := \Psi_{\text{LTL}} \wedge \Psi_P \wedge \Psi_N \wedge \Psi_D$ 
16: end while
17: return  $\varphi$ 

```

absence	existence	universality
$\mathbf{G}(\neg a_0)$	$\mathbf{F} a_0$	$\mathbf{G} a_0$
$\mathbf{G}(a_1 \rightarrow \mathbf{G}(\neg a_0))$	$(\mathbf{G}(\neg a_0)) \vee (\mathbf{F}(a_0 \wedge (\mathbf{F} a_1)))$	$\mathbf{G}(a_1 \rightarrow \mathbf{G} a_0)$
$\mathbf{F} a_1 \rightarrow (\neg a_0 \mathbf{U} a_1)$	$\mathbf{G}(a_0 \wedge ((\neg a_1) \rightarrow ((\neg a_1) \mathbf{U}(a_2 \wedge (\neg a_1))))))$	$\mathbf{F} a_1 \rightarrow (a_0 \mathbf{U} a_1)$
disjunction of common patterns		
$(\mathbf{F} a_2) \vee ((\mathbf{F} a_0) \vee (\mathbf{F} a_1))$		
$((\mathbf{G}(\neg a_0)) \vee (\mathbf{F}(a_0 \wedge (\mathbf{F} a_1)))) \vee ((\mathbf{G}(\neg a_3)) \vee (\mathbf{F}(a_2 \wedge (\mathbf{F} a_3))))$		
$(\mathbf{G}(a_0 \wedge ((\neg a_1) \rightarrow ((\neg a_1) \mathbf{U}(a_2 \wedge (\neg a_1)))))) \vee (\mathbf{G}(a_3 \wedge ((\neg a_4) \rightarrow ((\neg a_4) \mathbf{U}(a_5 \wedge (\neg a_4))))))$		

Table 1: Common LTL patterns used for generation of words.