

# Regex Decision Procedures in Extended RE#



Ian Erik Varatalu<sup>1</sup>[0000-0003-1267-2712]  
Margus Veanes<sup>2</sup>[0009-0008-8427-7977]  
Ekaterina Zhuchko<sup>1</sup>[0009-0004-8818-5042]  
Juhan Ernits<sup>1</sup>[0000-0002-4591-0425]



<sup>1</sup> Tallinn University of Technology, Tallinn, Estonia

<sup>2</sup> Microsoft Research, Redmond, USA

{ian.varatalu,ekzhuc,juhan.ernits}@taltech.ee, margus@microsoft.com

**Abstract.** We develop decision procedures for *extended regular expressions* in the new **ERE#** framework that uses *span semantics*, utilizing the power of *symbolic derivatives*. We prove a normal form theorem in Lean for **ERE#** that is closed under all Boolean operations and provides the basis for the given decision procedures. The tool is evaluated on existing SMT benchmarks for regexes that shows it to be the fastest solver to date – *often orders of magnitude faster than state-of-the-art* – albeit specialized for the *single-variable fragment* of string theory.

## 1 Introduction

The new and currently *fastest nonbacktracking* regex matcher **RE#** [46] supports match search with regexes extended with *intersection* ( $\&$ ), *complement* ( $\sim$ ), and restricted *lookarounds*. Extended operators in **RE#** provide a *separation of concerns* making it possible to express typical search patterns more naturally and compactly, as discussed at length in [46]. In this work, we develop decision procedures for the *Boolean closure* **ERE#** of **RE#** that enables various practical applications for *debugging*, *compiler optimizations*, and *verification* tasks.

Support for  $\&$  and  $\sim$  implies a need to decide *nonemptiness* of regexes as a basic *debugging* feature. For example, the regex  $.\ast\d.\ast\&\sim(..\ast\w.\ast)$  must match at least one digit but may not contain any word-letters – thus, it *matches nothing* because all digits are word-letters. Character classes can also be easily misunderstood due to subtle semantic differences between platforms or feature interactions in combination with common regex options, such as *case insensitivity* (typically inlined with  $(?i:..)$ ) where *equivalence* checking is essential for ensuring correctness and consistency of the intended behavior. It is possible that  $R_1 \equiv R_2$  but  $(?i:R_1) \not\equiv (?i:R_2)$ , e.g., for  $R_1 = [\wedge D]$  and  $R_2 = [\u0000-CE-\uFFFF]$ .

During derivative based compilation of regexes into DFAs, it is common to encounter unions of the form  $R_1 | R_2$  as target states of transitions. In order to reduce memory footprint, it is crucial to decide if  $R_i$  *subsumes*  $R_j$ , and if so, to reuse  $R_i$  as the sole target state, instead of introducing a new state for the union. Several basic rewrites necessary were already pointed out in [40], but *bounded loops* in particular are highly problematic [44]. E.g., the derivative of

a regex  $R = .*a\{k\}$  for  $a$  is  $R|a\{k-1\}$  that can be simplified to  $R$  because  $R$  subsumes  $a\{k-1\}$ , where  $a?$  abbreviates  $a|\varepsilon$ . Omitting such rewrites can quickly lead to a state space explosion up to size  $2^k$  – *observe that  $k = O(2^{|R|})$  here!* For DFA generation, such subsumption checking must be fast in practice.

Dealing with differences between PCRE and POSIX semantics for a regex is a common problem.<sup>3</sup> In PCRE, also known as *backtracking* or *leftmost-eager* semantics, the union operator  $|$  is *noncommutative* where in a regex  $R_1|R_2$ , a match for  $R_2$  is only sought when  $R_1$  fails to match. Note that this difference is relevant for *span* semantics that has recently been formalized in Lean [50], but irrelevant for *language* semantics (`IsMatch`) that is identical under both PCRE and POSIX. A technique to decide if a union  $R_1|R_2$  in PCRE has the same span semantics under POSIX, where union is *commutative*, is to decide match equivalence between  $R_1|R_2$  and  $R_1|R_2\&\sim(R_1\cdot_*)$ .

If the span semantics of a regex *differs* between PCRE and POSIX then the regex may contain unreachable cases under PCRE. E.g., the `BurntSushi/rebar` benchmarking tool [24] – *widely used for industrial PCRE matchers* – uses a dictionary benchmark containing unions such as `may|mayo`. The regex `may|mayo` will never match "mayo" for any input under PCRE semantics and has the exact same behavior as the regex `may`. Most industrial regex matchers use PCRE semantics, resulting in different behavior to what is intuitively assumed, as  $|$  is not union of languages in PCRE. It is a common programming error to define a regex with unintended behavior this way. For example, by using the above technique, `may|mayo&\sim(mayo_*)` effectively deletes the alternative `mayo` from `may|mayo` while `mayo|may` would remain intact because `may&\sim(mayo_*)`  $\equiv$  `may`.

The logic **ERE#**, that is introduced in Section 4, is a novel contribution and fundamental for many decision problems. In particular, it lifts **RE#** to the status of an *Effective Boolean Algebra* over spans. **RE#** is closed only under *intersection* and each regex admits a *linear* translation to a core regex of the form  $(?<=R_1)R_2(?=R_3)$  where no  $R_i$  contains lookarounds. The *span semantics*  $(u_1, u_2, u_3) \models (?<=R_1)R_2(?=R_3)$  intuitively means that  $u_1u_2u_3$  is a word such that  $u_i \in \mathcal{L}(R_i)$  where  $u_2$  is the main match with  $u_1$  and  $u_3$  as the surrounding context. For many decision problems, like subsumption, it became necessary to support regexes like  $(?<=R_1)R_2(?=R_3)\&\sim((?<=R'_1)R'_2(?=R'_3))$  that fall outside the fragment for matching in **RE#**.

Currently, **ERE#** semantics is not directly expressible in SMT-LIB as it would require support for lookarounds and span semantics. Section 7 proposes an *SMT-LIB format* for extending `RegLan` with *lookarounds* and *span semantics*.

**Contributions.** We introduce an extension **ERE#** of the **RE#** class that is Boolean closed and formalize a normal form theorem (Theorem 2) for it in *Lean* allowing us to develop decision procedures for **ERE#**, including *emptiness*, *subsumption*, and *equivalence*, that are of primary interest in the **ERE#** framework, as discussed above. It is also possible to use **ERE#** for pre-processing in SMT solvers. One can invoke **ERE#** from the simplifier in Z3 [39], to pre-process

<sup>3</sup> See for example <https://stackoverflow.com/questions/4733416>.

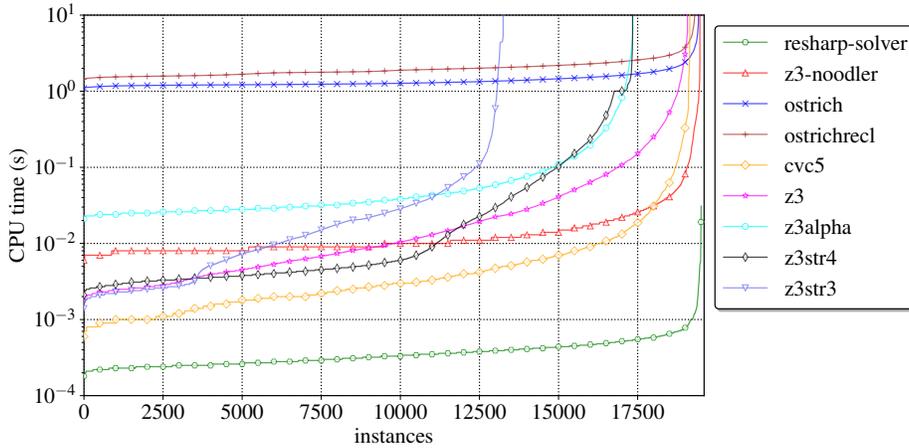


Fig. 1: Cactus plot of CPU time for 9 solvers and 19 509 SMT benchmarks. The y-axis is in **log-scale**. For **ERE#-solver** (**resharp-solver**) no benchmark needed  $>0.1$ s and less than 100 benchmarks needed  $>0.01$ s.

**RegLan** constraints, which can moreover be supported incrementally [12]. The main decision procedures for **ERE#** reduce, via the normal form theorem and the nonemptiness algorithm (Theorem 1), to nonemptiness in **ERE**.

We have conducted extensive evaluation on existing SMT benchmarks for the fragment of **ERE#** without lookarounds. The evaluation focuses on available SMT benchmarks that include all we could find on **RegLan** and benchmarks that reduce to the single string variable fragment of string constraints. The experiments demonstrate that **ERE#** consistently outperforms all state-of-the-art solvers, often by orders of magnitude, see Figure 1. The main reasons are: 1) *specialization to regex decision problems* combined with alphabet compression to bit-vectors; 2) aggressive rewrite rules that rely on symbolic derivatives and transition regex rewrites [43] – in particular lazy propagation rules for *intersection* and *complement*; 3) use of reversal theorem for **ERE<sub>≤</sub>** [50] – to decide nonemptiness of the *reverse* of a regex; 4) use of XOR as a Boolean operator directly supported by derivatives, for deciding equivalence in some cases.

**Artifact.** The paper is accompanied by an *artifact* [47] for Lean formalization of Theorem 2, and evaluation confirming the results in Figure 1 and Table 1.

## 2 Related Work

The development of string constraint solvers has a long-standing history, with prominent general-purpose SMT solvers such as Z3 [39] and CVC5 [5], alongside specialized solvers [10,19,36,11]. String constraint solving approaches vary widely, with methods based on a mixture of automata, word equations, and other

techniques. Many existing tools handle a broader range of constraints, including those beyond regular languages, such as context-free constraints. However, our work is not intended to compete with these general-purpose solvers; instead, we focus on a specialized fragment relevant to our setting. For a comprehensive overview of existing techniques and tools, we refer to the recent surveys [3,26]. Below we outline the key differences between the existing tools and our approach.

The solvers which support derivative-based lazy exploration of regular expressions are the sequence solver [43] in Z3 and CVC5 [5]. The former was the first to use transition regexes as part of the solver. However, our method leverages systematic simplifications and rewrites, enabling us to solve many SMT-LIB benchmarks in a fraction of a second, whereas the sequence solver fails to do so even after hours. Additionally, our solution is specialized for regexes, avoiding the overhead of supporting other theories. CVC5 [5] also uses a derivatives based approach [32] to solve regular expression constraints in its string solver [31,41], similarly relying on aggressive simplifications.

Kepler<sub>22</sub>, introduced in [30], employs a two-phase approach for solving string constraints. In the first phase, it constructs a cyclic reduction tree that represents all solutions to a conjunction of word equations and regular membership predicates. The second phase uses specialized procedures to infer length constraints from the set of all solutions, constituting a decision procedure for the straight line and quadratic fragments of string equations. While the authors report promising experimental results, we were not able to find an implementation to include in our evaluation. The tool NFA2SAT, presented in [33], performs an eager reduction of string operations and regular expression constraints into SAT, enabling incremental SAT solving. However, we were not able to find an artifact to include in our evaluation.

Several extensions to Z3 incorporate support for string constraints. The analysis in [9] focuses on fragments without string equality (i.e., without word equations), which aligns with the scope of our solver. Z3-Noodler [19] utilizes the MATA library [20] that implements basic algorithms for automata and also fast simulation reduction and antichain-based language inclusion checking. It uses novel techniques such as a stabilization-based procedure [18,14] for string-constraint solving. Z3str4 [36] builds on Z3str3 [10] and extends it with a string-to-bit-vector reduction. Z3alpha [34] is built on top of Z3 and the novelty lies in utilizing Monte Carlo Tree Search based SMT strategy synthesis. It ranked second in the 24-second performance category at SMT-COMP 2024, solving more tasks than all other solvers except Z3-Noodler.

Another notable tool is Ostrich [17], a solver that supports advanced regular expression features such as capture groups, lazy quantifiers, and anchors. Several extensions of Ostrich have been developed. For instance, [16] introduces a model based on prioritized streaming string transducers, while [27,28] builds on cost-enriched finite-state automata. Another recent extension to Ostrich is the tool SECO [29] based on parametric symbolic automata which extends symbolic automata to allow free variables on the transition guards. Ostrich has also been extended using Parikh images [23] to reason about word lengths. While

these Parikh image-based approaches have also been applied to accelerate string constraint solving [25], they fall outside the scope of this paper.

Equivalence algorithms for classical regular expressions have been studied in [1,2] based on *partial derivatives* [4]. The main algorithm was also implemented in [37] using Coq, and it operates incrementally by avoiding the construction of the full automaton during the process. The additional challenge in the case of **ERE#** is twofold: to support intersection and complement incrementally, and to work modulo large (or infinite) alphabets  $\mathcal{A}$  (given as an EBA), which we are currently investigating. Moreover, in the case of **ERE#**, equivalence should ideally also support lookarounds and therefore be based on span semantics.

### 3 Preliminaries

This section includes the main background notations and definitions used in the rest of the paper. To support large alphabets such as Unicode, we use *effective Boolean algebras* [22] (EBAs) to represent *character classes* symbolically via predicates, typical examples include: *digits*  $\backslash\mathbf{d}$ , *word-letters*  $\backslash\mathbf{w}$ , and *white-space characters*  $\backslash\mathbf{s}$ . Extended regular expressions or regexes support *intersection* and *complement* as well as *lookarounds* in the full class **ERE<sub>≤</sub>** [46,50]. The latter requires their semantics to be based on *locations* and *spans* due to context conditions imposed by lookarounds. Regexes use nested if-then-else terms called *transition regexes* [43,50] to represent their *symbolic derivatives*. The main intuition is that the symbolic derivative  $\delta(R)$  of a regex  $R$  is a transition regex representing a *partial evaluation* for  $R$  of the Brzozowski derivative [15]  $D_a(R)$  for  $a \in \Sigma$ . The *evaluation* of  $\delta(R)$  for  $a \in \Sigma$ ,  $\delta(R)[a]$ , equals  $D_a(R)$ . The class **ERE#**, defined in Section 4, forms a *proper* subclass of **ERE<sub>≤</sub>** but its match semantics is also based on spans and coincides with the match semantics of the whole **ERE<sub>≤</sub>**. The rest of this section defines these notions formally.

Before the formal definitions below, we illustrate the notions using the regex  $R = \_*\backslash\mathbf{d}_*\&_*\backslash\mathbf{w}_*\&_*\backslash\mathbf{s}_*$  that matches any string containing a digit, a word-letter, and a white-space character, where  $\_$  matches *all* characters and  $\&$  is regex *intersection*. The symbolic derivative  $\delta(R)$  of  $R$  is the transition regex:

$$\delta(R) = \mathbf{ite}(\backslash\mathbf{d}, \_*\backslash\mathbf{s}_*, \mathbf{ite}(\backslash\mathbf{w}, \_*\backslash\mathbf{d}_*\&_*\backslash\mathbf{s}_*, \mathbf{ite}(\backslash\mathbf{s}, \_*\backslash\mathbf{d}_*\&_*\backslash\mathbf{w}_*, R)))$$

The transition regex is illustrated as a binary decision tree in Figure 2 whose *conditions* are the highlighted predicates  $\backslash\mathbf{d}$ ,  $\backslash\mathbf{w}$ , and  $\backslash\mathbf{s}$ . The *else-case* branches of  $\delta(R)$  are *dashed* and the *leaves* of  $\delta(R)$  are regexes (including  $R$  itself). E.g.,  $\delta(R)[\text{'a'}] = \_*\backslash\mathbf{d}_*\&_*\backslash\mathbf{s}_*$  and  $\delta(R)[\text{'#'}] = R$ . When computing  $\delta(R)$  the alphabet EBA was used to eliminate unreachable subterms in  $\delta(R)$  through *cleaning* [43]. Therefore,  $\delta(R)$  was simplified by using the facts that  $\backslash\mathbf{d}$  implies  $\backslash\mathbf{w}$  (since all digits are word-letters), while  $\backslash\mathbf{s}$  is disjoint from both  $\backslash\mathbf{w}$  and  $\backslash\mathbf{d}$ .

**Effective Boolean Algebras.** An *Effective Boolean Algebra (EBA)* over an element universe  $\Sigma$  is a tuple  $(\Sigma, \mathcal{A}, \mathbf{f}, \perp, \_ , \sqcup, \sqcap, \mathbf{c})$  where  $\mathcal{A}$  is a set of *predicates*

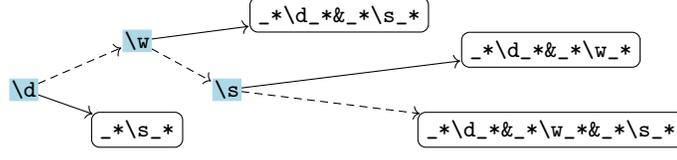


Fig. 2: Clean symbolic derivative of the regex  $_*\d_*&_*\w_*&_*\s_*$ .

that is closed under the Boolean connectives and contains  $\perp$  and  $\_$ . For  $a \in \Sigma$  and  $\alpha \in \mathcal{A}$  the *models* relation  $a \models \alpha$  with  $[[\alpha]] \stackrel{\text{DEF}}{=} \{a \in \Sigma \mid a \models \alpha\}$  obeys classical Tarski laws such that  $[[\perp]] = \emptyset$  and  $[[\_]] = \Sigma$ . For  $\alpha, \beta \in \mathcal{A}$  let  $\alpha \equiv \beta \stackrel{\text{DEF}}{=} [[\alpha]] = [[\beta]]$ . If  $\alpha \neq \perp$  then  $\alpha$  is *satisfiable* ( $\mathbf{SAT}(\alpha)$ ). All the connectives must be *computable* and  $\models$  must be *decidable*.

We write  $\mathcal{A}$  also for the EBA itself and say that  $\mathcal{A}$  is *decidable* if  $\mathbf{SAT}(\alpha)$  is decidable. We use  $(\Sigma, \mathcal{A}, \models, \perp, \_, \sqcup, \sqcap, \cdot^c)$  as a given *core alphabet* EBA. When working with regexes, it is a standard assumption that  $\cdot$  is a predicate denoting all characters *except* the newline character (predicate)  $\backslash n$ , i.e.,  $\_ \equiv \cdot \sqcup \backslash n$ .

**Locations and Spans.** A *location* is a pair of words in  $\Sigma^* \times \Sigma^*$ . Let  $w \in \Sigma^*$ . A *location in w* is a location  $(u, v)$  such that  $w = uv$ . The intuition is that a location in  $w$  specifies a border inside  $w$ . A location  $(\epsilon, w)$  is *initial* and a location  $(w, \epsilon)$  is *final*. For any location  $x$  let the *kind* of  $x$  be  $(IsInitial(x), IsFinal(x)) \in \mathbb{B} \times \mathbb{B}$  where  $\mathbb{B} = \{\mathbf{true}, \mathbf{false}\}$ .

A *span* is a triple of words in **Span**  $\stackrel{\text{DEF}}{=} \Sigma^* \times \Sigma^* \times \Sigma^*$ . A span  $\theta = (u, v, s)$  has two locations the *beginning* location  $beg(\theta) \stackrel{\text{DEF}}{=} (u, vs)$  and the *end* location  $end(\theta) \stackrel{\text{DEF}}{=} (uv, s)$ ;  $prefix(\theta) \stackrel{\text{DEF}}{=} u$  is the *prefix* of the span,  $suffix(\theta) \stackrel{\text{DEF}}{=} s$  is the *suffix* of the span, and  $match(\theta) \stackrel{\text{DEF}}{=} v$  is the *match* of the span. The *word* of the span is  $word(\theta) \stackrel{\text{DEF}}{=} uvs$ . We lift all the definitions to *sets* of spans as usual.

A *span in w* is a span  $(u, v, s)$  such that  $w = uvs$ . Intuitively, a span  $(u, v, s)$  represents a regex match in a word  $w$  with  $v$  as the primary matched substring of  $w$  where the prefix and the suffix are some sufficient lookahead context words. The *width* of a span  $(u, v, s)$  is the length  $|v|$  of its match  $v$ .

**Full Class ERE with Lookarounds.** The class  $\mathbf{ERE}_{\leq}$  is defined as follows. Members of  $\mathbf{ERE}_{\leq}$  are denoted here by  $R$ . *Concatenation*  $(\cdot)$  is often implicit by juxtaposition. All operators appear in order of precedence where *union*  $(\mid)$  binds weakest and *complement*  $(\sim)$  binds strongest. Let  $\psi \in \mathcal{A}$  and let  $m > 0$ .

$$R ::= \psi \mid \epsilon \mid R_1 \mid R_2 \mid R_1 \& R_2 \mid R_1 \cdot R_2 \mid R\{m\} \mid R^* \mid \sim R \mid \\ (?<=R) \mid (?<R) \mid (?=R) \mid (?!R)$$

The regex denoting *nothing* is the predicate  $\perp$ . Let  $R\{0\} \stackrel{\text{DEF}}{=} \epsilon$ . We write  $R+$  for  $R \cdot R^*$ . Union is also called *alternation* or *disjunction*. The regexes  $(?=R)$ ,  $(?!R)$ ,  $(?<=R)$ , and  $(?<R)$  are *lookarounds*;  $(?=R)$  is *(positive) lookahead*,

$(?!R)$  is *negative lookahead*,  $(?<=R)$  is (*positive*) *lookbehind*, and  $(?<!R)$  is *negative lookbehind*. Let  $\backslash A \stackrel{\text{DEF}}{=} (?<!_A)$  and  $\backslash z \stackrel{\text{DEF}}{=} (?!_z)$ .

**ERE** is the subclass of  $\mathbf{ERE}_{\leq}$  without lookarounds and  $\mathbf{ERE}_{\downarrow}$  extends **ERE** by allowing also the *start anchor*  $\backslash A$  and the *end anchor*  $\backslash z$  as primitive regexes.  $\mathbf{RE}_{\leq}$  (resp. **RE**) is the subclass of  $\mathbf{ERE}_{\leq}$  (resp. **ERE**) without  $\&$  and  $\sim$ .

**Match Semantics of ERE with Lookarounds.** The match semantics of regexes in  $\mathbf{ERE}_{\leq}$  uses spans [46,50]. Equivalent formulations of the semantics of **RE** with lookarounds appear originally in [38, Section 3.7] via derivation relations, and in [35,7] using a variant of spans. Let  $\theta = (u, v, s)$  be a span in a word  $w$ . Then  $\theta$  *models*  $R$  or  $R$  *matches*  $\theta$  is denoted by  $\theta \models R$ . We say that  $R$  *matches*  $w$  iff  $R$  matches some span in  $w$ . Recall that  $R\{0\} \stackrel{\text{DEF}}{=} \varepsilon$  and let  $m > 0$ .

$$\begin{aligned}
\theta \models L \mid R &\stackrel{\text{DEF}}{=} \theta \models L \vee \theta \models R \\
\theta \models L \& R &\stackrel{\text{DEF}}{=} \theta \models L \wedge \theta \models R \\
\theta \models \sim R &\stackrel{\text{DEF}}{=} \theta \not\models R \\
\theta \models R^* &\stackrel{\text{DEF}}{=} \exists n \geq 0 : \theta \models R\{n\} \\
(u, v, s) \models \varepsilon &\stackrel{\text{DEF}}{=} v = \varepsilon \\
(u, v, s) \models \psi &\stackrel{\text{DEF}}{=} |v| = 1 \wedge v \vDash \psi \\
(u, v, s) \models L \cdot R &\stackrel{\text{DEF}}{=} \exists x, y : v = xy \wedge (u, x, ys) \models L \wedge (ux, y, s) \models R \\
(u, v, s) \models R\{m\} &\stackrel{\text{DEF}}{=} \exists x, y : v = xy \wedge (u, x, ys) \models R \wedge (ux, y, s) \models R\{m-1\} \\
(u, v, s) \models (?<=R) &\stackrel{\text{DEF}}{=} v = \varepsilon \wedge (\varepsilon, u, s) \models \_ \cdot \_ \cdot R \\
(u, v, s) \models (?<!R) &\stackrel{\text{DEF}}{=} v = \varepsilon \wedge (\varepsilon, u, s) \not\models \_ \cdot \_ \cdot R \\
(u, v, s) \models (?=R) &\stackrel{\text{DEF}}{=} v = \varepsilon \wedge (u, s, \varepsilon) \models R \cdot \_ \cdot \_ \\
(u, v, s) \models (?!R) &\stackrel{\text{DEF}}{=} v = \varepsilon \wedge (u, s, \varepsilon) \not\models R \cdot \_ \cdot \_ \\
\mathcal{M}(R) &\stackrel{\text{DEF}}{=} \{\theta \in \mathbf{Span} \mid \theta \models R\} \\
L \equiv R &\stackrel{\text{DEF}}{=} \mathcal{M}(L) = \mathcal{M}(R)
\end{aligned}$$

Intuitively,  $(u, \varepsilon, s) \models (?=R)$  means that there exists a match of  $R$  *starting* from the location  $(u, s)$ , and  $(u, \varepsilon, s) \models (?<=R)$  means that there exists a match of  $R$  *ending* in the location  $(u, s)$ . For all lookarounds, the matched span has always 0 width, i.e., lookarounds are a generalized form of anchors.

Observe that  $\mathcal{M}(\_ \cdot \_) = \mathbf{Span}$ ,  $\mathcal{M}(\perp) = \emptyset$ , and all the Boolean connectives satisfy the EBA conditions of  $\mathbf{ERE}_{\leq}$ . Thus,  $(\mathbf{Span}, \mathbf{ERE}_{\leq}, \models, \perp, \_ \cdot \_, \mid, \&, \sim)$  is an EBA, since  $sp \models R$  is decidable for all  $sp \in \mathbf{Span}$  and  $R \in \mathbf{ERE}_{\leq}$  because  $\vDash$  is decidable in  $\mathcal{A}$ .

**Transition Regexes.** A *transition regex* is either a leaf  $R \in \mathbf{ERE}_{\downarrow}$ , or an ITE expression  $\mathbf{ite}(\psi, f, g)$  with *condition*  $\psi \in \mathcal{A}$ , *then-case*  $f$  and an *else-case*  $g$  that are transition regexes. The *evaluation* of  $f$  for  $a \in \Sigma$ , denoted by  $f[a]$ , is the leaf regex reached by  $a$ .

$$R[a] \stackrel{\text{DEF}}{=} R \quad \mathbf{ite}(\psi, f, g)[a] \stackrel{\text{DEF}}{=} \begin{cases} f[a], & \text{if } a \vDash \psi; \\ g[a], & \text{otherwise.} \end{cases}$$

All binary operators  $\diamond$  over regexes are lifted to transition regexes by propagating the operations into the leaves. Unary operators such as  $\sim$  and  $(?=)$  are

propagated analogously. Here  $R \in \mathbf{ERE}_{\perp}$ .

$$R \diamond \mathbf{ite}(\psi, f, g) \stackrel{\text{DEF}}{=} \mathbf{ite}(\psi, R \diamond f, R \diamond g) \quad \mathbf{ite}(\psi, f, g) \diamond h \stackrel{\text{DEF}}{=} \mathbf{ite}(\psi, f \diamond h, g \diamond h)$$

The set of leaves of a transition regex  $f$ , denoted by  $Lvs(f)$ , are the *reachable* leaves of  $f$ . In practice, transition regexes are kept *clean* [43] by construction, but here we filter out unreachable leaves explicitly. So  $\mathcal{A}$  is assumed *decidable*.

$$Lvs(f) \stackrel{\text{DEF}}{=} Lvs(\_, f) \quad Lvs(\psi, R) \stackrel{\text{DEF}}{=} \begin{cases} \{R\}, & \mathbf{if SAT}(\psi); \\ \emptyset, & \mathbf{otherwise.} \end{cases}$$

$$Lvs(\psi, \mathbf{ite}(\alpha, f, g)) \stackrel{\text{DEF}}{=} Lvs(\psi \sqcap \alpha, f) \cup Lvs(\psi \sqcap \alpha^c, g)$$

**Symbolic Derivatives in ERE with Anchors.** Here we define *nullability* and *symbolic derivatives* for regexes  $R \in \mathbf{ERE}_{\perp}$ . A (symbolic) derivative of  $R$  is either taken with respect to an *initial* (and nonfinal) location or a *noninitial* (and nonfinal) location, *final* locations are not used for computing derivatives. Nullability is defined for all locations.

We let the *kind*  $\kappa$  of a location  $x$  be  $(IsInitial(x), IsFinal(x)) \in \mathbb{B} \times \mathbb{B}$ . The three main kinds are  $\text{INI} = (\mathbf{true}, \mathbf{false})$ ,  $\text{MID} = (\mathbf{false}, \mathbf{false})$ , and  $\text{FIN} = (\mathbf{false}, \mathbf{true})$ . Using the definitions below we let  $Null(R) \stackrel{\text{DEF}}{=} Null_{\text{MID}}(R)$  and  $\delta(R) \stackrel{\text{DEF}}{=} \delta_{\text{MID}}(R)$ .

*Nullability.* We define nullability  $Null_{\kappa}(R)$  of a regex  $R \in \mathbf{ERE}_{\perp}$  relative to a *location kind*  $\kappa \in \mathbb{B} \times \mathbb{B}$ . Let  $m > 0$  and recall that  $R\{0\} \stackrel{\text{DEF}}{=} \varepsilon$ .

$$\begin{array}{ll} Null_{(initial, final)}(\backslash \mathbf{A}) \stackrel{\text{DEF}}{=} \mathbf{initial} & Null_{\kappa}(R \mid S) \stackrel{\text{DEF}}{=} Null_{\kappa}(R) \vee Null_{\kappa}(S) \\ Null_{(initial, final)}(\backslash \mathbf{z}) \stackrel{\text{DEF}}{=} \mathbf{final} & Null_{\kappa}(R \& S) \stackrel{\text{DEF}}{=} Null_{\kappa}(R) \wedge Null_{\kappa}(S) \\ Null_{\kappa}(\varepsilon) \stackrel{\text{DEF}}{=} \mathbf{true} & Null_{\kappa}(R \cdot S) \stackrel{\text{DEF}}{=} Null_{\kappa}(R) \wedge Null_{\kappa}(S) \\ Null_{\kappa}(R^*) \stackrel{\text{DEF}}{=} \mathbf{true} & Null_{\kappa}(R\{m\}) \stackrel{\text{DEF}}{=} Null_{\kappa}(R) \\ Null_{\kappa}(\psi) \stackrel{\text{DEF}}{=} \mathbf{false} & Null_{\kappa}(\sim R) \stackrel{\text{DEF}}{=} \neg Null_{\kappa}(R) \end{array}$$

So  $Null_{\text{INI}}(R)$  is applied in locations of kind INI and  $Null_{\text{FIN}}(R)$  is applied in locations of kind FIN. The trivial special case of the single empty location  $(\varepsilon, \varepsilon)$  in an empty word  $\varepsilon$  that is both initial and final is omitted from discussions.

*Symbolic Derivatives.* Symbolic derivatives, represented by transition regexes, are only evaluated for *nonfinal* locations, i.e.,  $\kappa \in \{\text{INI}, \text{MID}\}$  below. We apply the definition of derivatives from [48] and simultaneously lift the definition to transition regexes, similar to [43]. Let  $\psi \in \mathcal{A}$ ,  $m > 0$ , and  $\perp \in \{\backslash \mathbf{A}, \backslash \mathbf{z}\}$ .

$$\begin{array}{ll} \delta_{\kappa}(\perp) \stackrel{\text{DEF}}{=} \perp & \delta_{\kappa}(R^*) \stackrel{\text{DEF}}{=} \delta_{\kappa}(R) \cdot R^* \\ \delta_{\kappa}(\varepsilon) \stackrel{\text{DEF}}{=} \perp & \delta_{\kappa}(R\{m\}) \stackrel{\text{DEF}}{=} \delta_{\kappa}(R) \cdot R\{m-1\} \\ \delta_{\kappa}(R \& S) \stackrel{\text{DEF}}{=} \delta_{\kappa}(R) \& \delta_{\kappa}(S) & \delta_{\kappa}(\psi) \stackrel{\text{DEF}}{=} \mathbf{ite}(\psi, \varepsilon, \perp) \\ \delta_{\kappa}(R \mid S) \stackrel{\text{DEF}}{=} \delta_{\kappa}(R) \mid \delta_{\kappa}(S) & \delta_{\kappa}(R \cdot S) \stackrel{\text{DEF}}{=} \begin{cases} \delta_{\kappa}(R) \cdot S \mid \delta_{\kappa}(S), & \mathbf{if } Null_{\kappa}(R); \\ \delta_{\kappa}(R) \cdot S, & \mathbf{otherwise.} \end{cases} \\ \delta_{\kappa}(\sim R) \stackrel{\text{DEF}}{=} \sim \delta_{\kappa}(R) & \end{array}$$

Observe that, by definition, all regex operators are lifted to transition regexes.

## 4 Decision Procedures for $\mathbf{ERE\#}$

The focus of the paper is on the subclass  $\mathbf{ERE\#}$  of  $\mathbf{ERE}_{\leq}$  that contains  $\mathbf{ERE}_{\downarrow}$  combined with a restricted fragment of lookarounds. In the formal definition below  $L$  defines  $\mathbf{ERE}_{\downarrow}$ , and  $R$  defines  $\mathbf{ERE\#}$ . Let  $\psi \in \mathcal{A}$  and  $m > 0$ .

$$\begin{aligned} L &::= \psi \mid \varepsilon \mid \backslash \mathbf{A} \mid \backslash \mathbf{z} \mid L_1 \mid L_2 \mid L_1 \& L_2 \mid L_1 \cdot L_2 \mid L\{m\} \mid L^* \mid \sim L \\ R &::= L \mid (?<=L) \cdot R \mid (?<L) \cdot R \mid R \cdot (?=L) \mid R \cdot (?!L) \mid R_1 \mid R_2 \mid R_1 \& R_2 \mid \sim R \end{aligned}$$

We show below that all regexes in  $\mathbf{ERE\#}$  have a normal form that is a union of *core regexes* that are regexes of the form  $(?<=L_1) \cdot L_2 \cdot (?=L_3)$  where  $L_i \in \mathbf{ERE}_{\downarrow}$ , with  $L_2$ ,  $(?<=L_1) \cdot L_2$ ,  $L_2 \cdot (?=L_3)$  as special cases because  $(?<=\varepsilon) \equiv (?=\varepsilon) \equiv \varepsilon$ . The definition of  $\mathbf{ERE\#}$  properly subsumes the definition of  $\mathbf{RE\#}$  in [46] where  $\mathbf{RE\#}$  contains all core regexes and is only closed under  $\&$ .

The match semantics of  $\mathbf{ERE\#}$  is based on  $\mathbf{ERE}_{\leq}$ . In particular, it follows that  $(\mathbf{Span}, \mathbf{ERE\#}, \models, \perp, \_*, \mid, \&, \sim)$  is an EBA. The main decision procedures we are focusing on for  $\mathbf{ERE\#}$  are *emptiness*, *subsumption*, and *equivalence*.

### 4.1 Deciding Nonemptiness of Core Regexes in $\mathbf{ERE\#}$

Here we consider nonemptiness of *core regexes* in  $\mathbf{ERE\#}$ . We later show how the general case of nonemptiness of all regexes in  $\mathbf{ERE\#}$  reduces to nonemptiness of core regexes by showing that  $\mathbf{ERE\#}$  is a *decidable* EBA. The nonemptiness algorithm builds on symbolic derivatives and transition terms and can be abstractly formulated as a fixpoint procedure that relies on the associativity, commutativity, and idempotence (ACI) of regex union, which guarantees finiteness of the state space and thus termination.

For a union regex let  $Set(R_1 \mid R_2) \stackrel{\text{DEF}}{=} Set(R_1) \cup Set(R_2)$ . Let  $Set(\perp) \stackrel{\text{DEF}}{=} \emptyset$  and for all other regexes  $R$  let  $Set(R) \stackrel{\text{DEF}}{=} \{R\}$ . Let  $f$  be a transition regex. The set of all *states* of  $f$  is the set of all  $q \in Set(\ell)$  for  $\ell \in Lvs(f)$ , i.e.,

$$States(f) \stackrel{\text{DEF}}{=} \bigcup_{\ell \in Lvs(f)} Set(\ell)$$

For all  $R = (?<=L_1) \cdot L_2 \cdot (?=L_3) \in \mathbf{ERE\#}$ , where  $L_i \in \mathbf{ERE}_{\downarrow}$ , we decide nonemptiness of  $R$  by reducing it to nonemptiness in  $\mathbf{ERE}_{\downarrow}$  as follows.

$$IsNonempty((?<=L_1) \cdot L_2 \cdot (?=L_3)) \stackrel{\text{DEF}}{=} IsNonempty(\_ * L_1 \cdot L_2 \cdot L_3 \cdot \_ *)$$

Let  $L \in \mathbf{ERE}_{\downarrow}$ . The function  $IsNonempty(L)$ , unless  $L$  is trivially nullable, computes reachable states from  $L$  and returns **true** upon reaching a nullable state.

$IsNonempty(L) \stackrel{\text{DEF}}{=} \mathbf{if} \ Null_{(\mathbf{true}, \mathbf{true})}(L) \ \vee \ Null_{\text{INI}}(L) \ \mathbf{return} \ \mathbf{true}$   
 $\rho(L) \leftarrow \delta_{\text{INI}}(L); \ Q \leftarrow \{L\} \cup States(\rho(L))$   
 $\mathbf{while} \ \nexists q \in Q : (Null(q) \ \vee \ Null_{\text{FIN}}(q)) \ \wedge \ \exists q \in Q \setminus Dom(\rho) \ \mathbf{do}$   
 $\rho(q) \leftarrow \delta(q); \ Q \leftarrow Q \cup States(\rho(q))$   
 $\mathbf{return} \ \exists q \in Q : Null(q) \ \vee \ Null_{\text{FIN}}(q)$

$$\begin{array}{c}
\text{NOT-ELIM} \frac{\sim((?<=X)Y(?=Z))}{(?<!X)_* | \sim Y | \sim (?!Z)} \quad \text{L-DISTR} \frac{(Y|Z)X}{YX|ZX} \quad \text{R-DISTR} \frac{X(Y|Z)}{XY|XZ} \\
\text{NLA-ELIM} \frac{(?!X)}{(?=\sim(X_*)\backslash z)} \quad \text{LA-JOIN} \frac{(?=X)(?=Y)}{(?=X_*\&Y_*)} \quad \text{LB-JOIN} \frac{(?<=X)(?<=Y)}{(?<=_*X\&_*Y)} \\
\text{NLB-ELIM} \frac{(?<!X)}{(?<=\backslash A^{\sim}(_*X))} \quad \text{AND-ELIM} \frac{(?<=X)Y(?=Z)\&(?!X')Y'(?=Z')}{(?<=X)(?!X')(Y\&Y')(?!Z)(?!Z')}
\end{array}$$

Fig. 3: Main inference rules in  $\mathbf{ERE}_{\leq}$  used in Theorem 2. Let  $\text{NOT-ELIM\#} = \text{NOT-ELIM} \circ \text{NLB-ELIM} \circ \text{NLA-ELIM}$ , and  $\text{AND-ELIM\#} = \text{AND-ELIM} \circ \text{LB-JOIN} \circ \text{LA-JOIN}$ .

*Example 1.* Let  $R := (?<=\backslash A)\varepsilon(?!z)$ . Then  $\text{Null}_{\text{INI}}(\sim \backslash A^{\sim} \backslash z_*) = \mathbf{true}$ . Here  $R$  only matches spans of 0 width whose start location is initial and end location in nonfinal, i.e.,  $\mathcal{M}(R) = \{(\varepsilon, \varepsilon, v) \mid v \in \Sigma^+\}$ .

Let  $R := (?<=\backslash A)_*(?!z)$ . Then  $L = \sim \backslash A \backslash z_*$  and  $\rho(L) = L | \backslash z_*$ . So initially  $Q = \{L, \backslash z_*\}$  because  $\text{States}(L | \backslash z_*) = \{L, \backslash z_*\}$  and we have that  $\text{Null}_{\text{FIN}}(\backslash z_*) = \mathbf{true}$ . Note that  $\mathcal{M}(R) = \{(\varepsilon, a, \varepsilon) \mid a \in \Sigma\}$ .  $\square$

**Theorem 1 (NonEmptiness)**  $\mathcal{M}(R) \neq \emptyset \Leftrightarrow \text{IsNonempty}(R)$  holds for all core regexes  $R$  in  $\mathbf{ERE\#}$ .

$\text{IsNonempty}(R)$  can be extended to produce a *witness* when  $\mathbf{true}$ . In most precise form, the produced witness can be in form of a *symbolic span*  $(\alpha_1, \alpha_2, \alpha_3) \in \mathcal{A}^* \times \mathcal{A}^* \times \mathcal{A}^*$  such that, for all spans  $(u_1, u_2, u_3)$ , such that  $u_i \vDash \alpha_i$ , it holds that  $(u_1, u_2, u_3) \vDash R$ , where  $(a_i)_{i<n} \vDash (\psi_i)_{i<m}$  stands for  $n = m$  and  $\bigwedge_{i<n} (a_i \vDash \psi_i)$ .

## 4.2 Lookaround Normal Form of $\mathbf{ERE\#}$

We start by providing a collection of equivalence preserving inference rules in  $\mathbf{ERE}_{\leq}$ , see Figure 3, that are used to rewrite regexes in  $\mathbf{ERE\#}$  into the desired normal form. All rules have been formalized and proved correct in Lean. Observe that all rules preserve  $\mathbf{ERE\#}$ , i.e., lookaround bodies remain in  $\mathbf{ERE}_{\downarrow}$ .

The key new rule is NOT-ELIM. The other rules involving lookarounds are derived from similar rewrites in [46] where the normal form is *linear* in the size of the original  $\mathbf{RE\#}$  formula. In the case of  $\mathbf{ERE\#}$  the lookaround normal form can be *exponential* in the size of the original formula.

Let  $R$  be a Boolean combination of *core regexes*. We define the *negation normal form*  $\text{NNF}(R)$  of  $R$  where regex complement  $\sim$  has been propagated via de Morgan extended with the rule NOT-ELIM#. It follows that no subregex of  $\text{NNF}(R)$  outside  $\mathbf{ERE}_{\downarrow}$  is negated and all lookarounds are positive.

Many further simplification laws are used as rewrites. The basic ones are treating  $(\sim, \perp)$  as the units of  $(\&, |)$  and as the zeros of  $(|, \&)$ , and  $\perp$  is also the zero of  $\cdot$  and  $\varepsilon$  is the unit of  $\cdot$ . Further rules include  $(?! \varepsilon) \equiv \perp$  and  $(?! \varepsilon) \equiv \perp$ , as well as  $(?= \perp) \equiv \perp$  and  $(?<= \perp) \equiv \perp$ . Also,  $\sim \backslash A \equiv \backslash A$  and  $\backslash z_* \equiv \backslash z$ .

Further practical rewrites for negative lookarounds with body  $\psi \in \mathcal{A}$  are  $(?! \psi) \equiv (?<= \psi^c | \backslash A)$  and  $(?! \psi) \equiv (?= \psi^c | \backslash z)$  as equivalent but simplified variants of NLB-ELIM and NLA-ELIM, respectively.

*Example 2.* The regex  $(?!abc\z)$  is intuitive, it states that the suffix of a match must not be  $abc$ . So  $(?=\sim(abc\z)\z)$  is equivalent but less intuitive.

To illustrate some rules, consider  $\sim((?<=a)_*b_*|_*c_*)$  in **ERE#**. Here we also make use of  $\sim(_*b_*) \equiv [\wedge b]^*$  and  $\sim(_*c_*) \equiv [\wedge c]^*$ , as well as  $[\wedge b]^*[\wedge c]^* \equiv [\wedge bc]^*$  (using standard negated character class notation).

$$\begin{aligned}
\sim((?<=a)_*b_*|_*c_*) &\stackrel{\text{DEMORGAN}}{\equiv} \sim((?<=a)_*b_*)\&\sim(_*c_*) \\
&\stackrel{\text{NOT-ELIM}}{\equiv} ((?<!a)_*|\sim(_*b_*))\&\sim(_*c_*) \\
&\stackrel{\text{NLB-ELIM}}{\equiv} ((?<=[\wedge a]|\wedge A)_*|\sim(_*b_*))\&\sim(_*c_*) \\
&\equiv ((?<=[\wedge a]|\wedge A)_*|[\wedge b]^*)\&[\wedge c]^* \\
&\equiv (?<!a)[\wedge c]^*|[\wedge bc]^*
\end{aligned}$$

For example,  $(\varepsilon, b, \varepsilon) \models (?<!a)[\wedge c]^*$  but  $(\varepsilon, b, \varepsilon) \not\models (?<=a)_*b_*|_*c_*$ .  $\square$

### 4.3 Formalization in Lean

We now present the key result that establishes the correctness of the Lookaround Normal Form (LNF). Our approach follows the formalization presented in [50], which defines the match semantics for the entire class of regular expressions **ERE<sub>≤</sub>**. In this work, we explicitly handle three distinct types of match semantics, each corresponding to different subclasses of regular expressions: **ERE<sub>↓</sub>**, **ERE#**, and **ERE<sub>≤</sub>**, and establish their relationships to ensure correct correspondence between their semantics. Additionally, we introduce an internal conversion between the **RESharp** type and its positive fragment, **PosRESharp**, which excludes complements and negative lookarounds.

First, we present some auxiliary definitions. The central definition in our formalization is the `lnf` function, which converts a regex in **ERE#** into a list of core regexes. The normalization function `LNF` in Theorem 2 is defined in terms of `lnf`. The `lnf` function is defined inductively and implements the inference rules presented in Figure 3. Components such as `la_join`, `lb_join`, and `intersection` closely follow the corresponding rules.

The `lnf` function definition, returning a list of core regexes, is as follows.<sup>4</sup>

```

def lnf (r : RESharp A) : List (CoreRegex A) :=
match r with
| Ere r           => [EREa_to_CoreRegex r]
| Lookahead r la => map (la_join la) (lnf r)
| Lookbehind lb r => map (lb_join lb) (lnf r)
| NLookahead r la => map (la_join (nLookahead la)) (lnf r)
| NLookbehind lb r => map (lb_join (nLookbehind lb)) (lnf r)
| Alt l r         => lnf l ++ lnf r
| Inter l r       => productWith intersection (lnf l) (lnf r)
| Compl r         => takeNegations (lnf r)

```

<sup>4</sup> The full Lean formalization is available in the supplemental material.

The most intricate case involves the complement operation. To handle this case more easily and to simplify the proof structure, we introduce an auxiliary function, `takeNegations`, that specifically deals with negated expressions.

```
def takeNegations (rs : List (CoreRegex  $\mathcal{A}$ )) : List (CoreRegex  $\mathcal{A}$ ) :=
match rs with
| []      => [EREA_to_CoreRegex _*]
| c :: cs => let cs' := takeNegations cs
  map (.lb_join (nLookbehind c.left)) cs' ++
  map (.la_join (nLookahead c.right)) cs' ++
  map (intersection (EREA_to_CoreRegex (~c.regex))) cs'
```

Intuitively, the three lists correspond to the inference rule NOT-ELIM where any of the three components of the union can be satisfied in order to satisfy the negation of the whole expression.

Let us give the remaining definition, LNF, in the Theorem 2.

```
def LNF (r : RESharp  $\mathcal{A}$ ) : RESharp  $\mathcal{A}$  := re_sum_pos (map sem (lnf r))
```

It first transforms the input regex into its Lookaround Normal Form (LNF) using the `lnf` function. Then, it applies the `sem` function to each core regex in the LNF list, converting them into their corresponding **ERE#** representation e.g., the core regex  $(L_1, L_2, L_3)$  is converted into  $(?<=L_1) \cdot L_2 \cdot (?=L_3)$ . Finally, `re_sum_pos` combines the resulting list of **ERE#** expressions with the union operator.

Finally, we are ready to state the main theorem where **RESHarp** denotes **ERE#**.

### Theorem 2 (LNF)

```
theorem lnf_correct {R : RESharp  $\mathcal{A}$ } {sp : Span  $\Sigma$ } : sp  $\models$  R  $\leftrightarrow$  sp  $\models$  LNF R
```

Theorem 2 states that the match semantics of a regex  $R$  remains unchanged when it is transformed into its equivalent Lookaround Normal Form (LNF). In other words, converting a regex into LNF preserves its matching behavior. The proof proceeds by induction on  $R$ , where we demonstrate that the normalization function behaves as expected for each case.

**Theorem 3 (Decidability)** *If  $\mathcal{A}$  is decidable then so is **ERE#** modulo  $\mathcal{A}$ .*

*Proof.* Let  $R \in \mathbf{ERE\#}$ . Then  $\mathcal{M}(R) \neq \emptyset \Leftrightarrow \exists S \in \text{LNF}(R) : \mathcal{M}(S) \neq \emptyset$  holds by Theorem 2 and  $\mathcal{M}(S) \neq \emptyset \Leftrightarrow \text{IsNonempty}(S)$  holds by Theorem 1, where  $\text{IsNonempty}(S)$  uses *cleaning* which requires  $\mathcal{A}$  to be decidable.  $\square$

Decidability of **RE** extended with arbitrary lookaheads has been shown in [8]. Note that **ERE#** supports a restricted form of lookaheads, e.g., nested lookaheads are not supported and **ERE#** is not closed under concatenation. On the other hand, **ERE#** supports also restricted lookbehinds and is closed under reversal. Although we have provided a decision procedure for **ERE#**, an efficient symbolic decision procedure for the full class **ERE<sub>≤</sub>** remains an open problem.

#### 4.4 Subsumption and Equivalence in ERE#

Let  $R, S \in \mathbf{ERE\#}$ . Then  $S$  *subsumes*  $R$  or  $R$  *is subsumed by*  $S$ , denoted by  $R \sqsubseteq S$  means that  $\mathcal{M}(R) \subseteq \mathcal{M}(S)$ . Subsumption is equivalent to *emptiness* of  $R \& \sim S$ , since  $\mathcal{M}(R \& \sim S) = \emptyset \Leftrightarrow \mathcal{M}(R) \setminus \mathcal{M}(S) = \emptyset \Leftrightarrow \mathcal{M}(R) \subseteq \mathcal{M}(S)$ .

One can also consider various *language semantics* of  $R$  such as  $\text{word}(\mathcal{M}(R))$  and  $\text{match}(\mathcal{M}(R))$  and related equivalence relations. In general, those are weaker than  $\equiv$ . E.g.,  $\mathcal{M}(\text{?<=}\backslash\mathbf{Aa}\backslash\mathbf{z}) = \{(\mathbf{a}, \varepsilon, \varepsilon)\}$  and  $\mathcal{M}(\text{?}=\backslash\mathbf{Aa}\backslash\mathbf{z}) = \{(\varepsilon, \varepsilon, \mathbf{a})\}$  while the *word* language is  $\{\mathbf{a}\}$  and the *match* language is  $\{\varepsilon\}$  in both cases. The important special case is for regexes  $\backslash\mathbf{AR}\backslash\mathbf{z}$  whose language  $\text{match}(\mathcal{M}(\backslash\mathbf{AR}\backslash\mathbf{z}))$  corresponds exactly to  $\mathcal{M}(\backslash\mathbf{AR}\backslash\mathbf{z})$ .

One method to decide equivalence  $R \equiv S$  is to decide that both  $R \sqsubseteq S$  and  $S \sqsubseteq R$  hold, that is emptiness of  $R \& \sim S \mid S \& \sim R$ . A different method is to use an additional Boolean operator for *symmetric difference* (XOR)  $\oplus$  whose derivative rule in  $\mathbf{ERE\#}$  would remain the same as for all the other binary operators:

$$\delta_\kappa(R \oplus S) = \delta_\kappa(R) \oplus \delta_\kappa(S) \quad \text{Null}_\kappa(R \oplus S) \stackrel{\text{DEF}}{=} (\text{Null}_\kappa(R) \neq \text{Null}_\kappa(S))$$

Thus,  $\text{IsNonempty}(R \oplus S) \Leftrightarrow R \neq S$ . While this works for  $R, S \in \mathbf{ERE\#}$ , it is unclear if this method generalizes to the whole  $\mathbf{ERE\#}$  in a “useful” manner, since  $\text{LNF}(R \oplus S)$  would have to break the problem into separate disjuncts. The  $\oplus$  operator also has a dual *XNOR* operator  $\odot$  so that  $\sim(R \oplus S) \equiv \sim R \odot \sim S$ :

$$\delta_\kappa(R \odot S) = \delta_\kappa(R) \odot \delta_\kappa(S) \quad \text{Null}_\kappa(R \odot S) \stackrel{\text{DEF}}{=} (\text{Null}_\kappa(R) = \text{Null}_\kappa(S))$$

So  $\oplus$  and  $\odot$  can be used like any other binary operator in any regex.

*Example 3.* Given the regex  $q_0 = \backslash\mathbf{A}([\mathbf{ab}] + \&\sim(\_ * \mathbf{aa} \_)) \backslash\mathbf{z}$  and the regex  $r_0 = \backslash\mathbf{A}(\mathbf{a}(\mathbf{b} + \mathbf{a}?) * \mid \mathbf{b}(\mathbf{a}?( \backslash\mathbf{z} \mid \mathbf{b}+)) *) \backslash\mathbf{z}$ , we consider their equivalence. Both regexes match all strings in  $[\mathbf{ab}]^+$  without  $\mathbf{aa}$  as a substring. While  $q_0$  is quite transparent in this regard,  $r_0$  is less clear but avoids  $\&$  and  $\sim$ . We prove their equivalence below, by using  $\oplus$ . We first describe their symbolic derivatives and resulting (symbolic) DFAs separately, as illustrated in Figure 4. For example, for the regex  $q_0$  we get that  $\delta_{\text{INI}}(q_0) = \mathbf{ite}(\mathbf{a}, q_2, \mathbf{ite}(\mathbf{b}, q_1, \perp)) = \delta(q_1)$ , and for the regex  $r_0$  we get that  $\delta_{\text{INI}}(r_0) = \mathbf{ite}(\mathbf{a}, r_1, \mathbf{ite}(\mathbf{b}, r_2, \perp))$  and  $\delta(r_1) = \mathbf{ite}(\mathbf{b}, r_3, \perp)$ , etc., where the leaf regexes are shown as states in Figures 4a and 4b. Now, returning to equivalence, the DFA for  $r_0 \oplus q_0$  looks similar to Figure 4b, e.g.,

$$\begin{aligned} \delta_{\text{INI}}(r_0 \oplus q_0) &= \mathbf{ite}(\mathbf{a}, q_2, \mathbf{ite}(\mathbf{b}, q_1, \perp)) \oplus \mathbf{ite}(\mathbf{a}, r_1, \mathbf{ite}(\mathbf{b}, r_2, \perp)) \\ &\equiv \mathbf{ite}(\mathbf{a}, q_2 \oplus r_1, \mathbf{ite}(\mathbf{b}, q_1 \oplus r_2, \perp)) \end{aligned}$$

except that the DFA has no accepting states. E.g.,  $\text{Null}_{\text{FIN}}(q_2 \oplus r_1) = \mathbf{false}$  because both  $q_2$  and  $r_1$  are nullable, similarly for all the other states. Thus  $r_0 \equiv q_0$ .  $\square$

An algorithm for regex equivalence based on  $\oplus$  would not need to construct a complete DFA for  $\oplus$ , in particular when a witness of inequivalence is not needed. Such an algorithm is ongoing work that falls outside the scope of this paper.

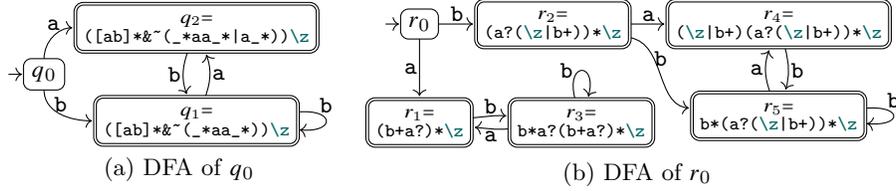


Fig. 4:  $q_0 = \backslash A([ab]^* \& \sim (_*aa_*) \backslash z$  and  $r_0 = \backslash A(a(b+a?)^* | b(a?( \backslash z | b+)) * \backslash z$ . The *sink* state  $\perp$  is implicit – from each state  $s$  there is a symbolic transition to the sink state for all the other characters, e.g.,  $\delta(q_1)[c] = \perp$  and  $\delta(r_1)[a] = \perp$ .

## 5 Implementation

Our implementation of the decision procedures is done in the Rust programming language. The implementation is available on GitHub [45]. It supports a subset of the Unicode Strings theory in SMT-LIB [13], with a parser that converts smt2 format into our internal representation of regexes.

All the supported operations are implemented as operations on **ERE#**, and we support all the **RegLan** operations in the SMT-LIB standard. Operations that can not be represented as regexes in a straightforward manner such as **str.replace** are considered out of scope, as well as operations on a higher level of abstraction, such as the loop  $R\{n\}$ , where  $n$  is not a constant. In the formal semantics, the regex  $R\{m, n\}$ , where  $0 < m \leq n$ , stands for  $R\{m\} \cdot (R|\epsilon)\{n-m\}$ .

The operations supported in the implementation include the following.

- All operations starting with **re.** on constants or other regexes, such as **re.++**, **re.union**, **re.\***, etc.
- **ERE** membership **str.in\_re** on variables and constants
- String assertions that translate intuitively into **ERE**, such as:
 

$\frac{\text{str.len } x < 5}{\_ \{0,4\}}$	$\frac{\text{str.prefixof "abc" } x}{abc\_ *}$	$\frac{\text{str.contains "abc" } x}{\_ *abc\_ *}$
--	--	--
- Boolean operations such as **and**, **or**, **not** on any of the above, e.g.,
 

$\frac{(\text{str.len } x > 3) \text{ and } (\text{str.len } x < 9)}{\_ \{4,\} \& \_ \{0,8\}}$	$(\equiv \_ \{4,8\})$
--	-----------------------
- Standard conversions and utilities such as **str.to\_int**, **str.to\_re**, **str.++**

### 5.1 Representation of Regexes

In the implementation, each regex is represented as a tree structure of nodes, where each node has:

1. a unique identifier (unsigned 32-bit integer)
2. a kind specifier (e.g., predicate, union, concatenation, etc.)
3. the unique identifiers of the left and right child (or 0 if missing)
4. various additional information (e.g., bit-flags, and relevant character set)

Data stored as additional information includes a compact integer encoding of a predicate  $\varphi_r$  for each node  $r$  that approximates the relevant characters in  $r$ , as defined in [46, Section 5.3] with related rewrites. Additionally,  $r$  maintains several bit-flags that can be inferred during construction, such as whether the language contains the empty string (i.e.,  $r$  is nullable), the presence of anchors in  $r$ , or whether  $r$  contains extended operations such as  $\&$  or  $\sim$ .

When working with UTF8 representation, the character EBA  $\mathcal{A}$  is over *bytes* as  $\Sigma$ . The representation of valid UTF8 strings is the regex  $R_{\text{UTF8}}$ :

$$([\text{x}00\text{-}\text{x}7\text{F}] | [\text{x}C0\text{-}\text{x}DF]\beta | [\text{x}E0\text{-}\text{x}EF]\beta\{2\} | [\text{x}F0\text{-}\text{x}F7]\beta\{3\})^*$$

where  $\beta = [\text{x}80\text{-}\text{x}BF]$  is the *continuation byte* predicate. When working with any regex  $R \in \mathbf{ERE}_{\perp}$ , the intersection  $R \ \& \ \backslash AR_{\text{UTF8}} \ \backslash z$  restricts the accepted word language to  $\mathcal{L}(R_{\text{UTF8}})$  that is a proper subset of  $\Sigma^*$ .

The predicates  $\psi$  in  $\mathcal{A}$  are represented as 256-bit bit-vectors, allowing for efficient constant-time Boolean operations. Each  $i \in \Sigma$  corresponds to the  $i$ 'th bit of  $\psi$ .  $\mathcal{A}$  has trivial checks for  $\psi \equiv \perp$  and for  $\psi \equiv \_$  since  $\perp = 0$  and  $\_ = 2^{256} - 1$ .

To accelerate the bit-vector operations further, we perform the operations in parallel over the 64-bit integers using Single Instruction Multiple Data (SIMD) instructions available in modern CPUs, e.g., *AVX2* or *SSE4.2*.

We assign a unique identifier to each predicate, allowing us to represent the predicate of a node as a 32-bit integer instead of a 256-bit bit-vector. This reduces the maximum number of unique predicates in a given regex to  $2^{32}$ , but in practice this is a reasonable trade-off, as it is very unlikely that a regex would contain more than a thousand unique character sets, let alone  $2^{32}$ .

The identifier of the node represents the syntactic structure of the regex, and is used to memoize the results of the operations on the regexes. The integers are assigned sequentially during construction, for example, the regex  $(\mathbf{a|b})$  would be assigned the integer 3, and point to the corresponding leaf nodes 1 and 2, representing the characters  $\mathbf{a}$  and  $\mathbf{b}$  respectively. To prevent duplicate identifiers, we use a hash map to store the regexes that have been constructed so far.

While not strictly relevant to the theory, we also have several rules for normalization, such as always having the left child of a union or intersection operation be the smaller identifier. Thus, the regexes  $(\mathbf{a|b})$  and  $(\mathbf{b|a})$  would be normalized to the same regex, whichever is constructed first. We also enforce that concatenations can only contain inner concatenations on the right to prevent ambiguous nesting. Together, these rules ensure that the syntactic construction of the regex is unique up to commutativity and associativity, but this does not prevent the construction of equivalent regexes in terms of the language they represent. For example,  $(\varepsilon|\mathbf{b})\{50\}$  and  $(\mathbf{b}\{0,25\})\{2\}$  represent the same language, but are constructed as different regexes – to detect equivalence between such examples, we would need to use the equivalence decision procedure.

## 5.2 Rewrite Rules and Simplifications

An important part of the implementation is the simplification and memoization of regexes during construction. The simplifications are applied bottom-up, i.e.,

$$\frac{R}{\perp}(\mathcal{M}(R)=\emptyset) \quad \frac{R}{\_*}(\mathcal{M}(R)=\mathbf{Span}) \quad \frac{R_1 \mid R_2}{R_2}(\mathcal{M}(R_1) \subseteq \mathcal{M}(R_2)) \quad \frac{R_1 \& R_2}{R_1}(\mathcal{M}(R_1) \subseteq \mathcal{M}(R_2))$$

Fig. 5: Main rewrite rules at the level of *match semantics*.

when a regex is constructed, the children of the regex are already simplified. The simplifications are based on the rules in [46], but in a more general way, as, in addition to the algebraic rules, we can perform language-level operations on the regexes directly.

For example, the regex  $\_ * a \_ * b \_ * \mid \_ * b \_ *$  would be simplified to  $\_ * b \_ *$  during construction, as the left branch of the union is subsumed by the right branch. In other words, the language of  $\_ * a \_ * b \_ *$  is a subset of the language of  $\_ * b \_ *$ . Subsumption plays a crucial role in the construction of intersections and unions, as it allows us to prevent construction of redundant definitions, often considerably reducing the size of the regex.

We can also check satisfiability of a regex during construction, which allows us to rewrite unsatisfiable regexes to  $\perp$ , which represents the empty language.

For the key rewrite rules on the *match semantics* level, which are used in the implementation, see Figure 5.

Performing these operations during construction gives us a very important property for optimization, which is that any remaining regex after construction is either  $\perp$  or satisfiable. This allows us to skip many unnecessary satisfiability checks during the decision procedure.

This satisfiability check (leftmost rule in Figure 5) only needs to be performed when constructing a concatenation with an anchor, an intersection, or complement, as these operations can result in an unsatisfiable regex, given that the children of the operation are satisfiable. The other operations will always result in a satisfiable regex, which we propagate as bit-flags in the implementation.

While the construction of a concatenation node with a satisfiable node and an anchor may result in an unsatisfiable regex, this satisfiability check takes constant time, as it only requires to check whether the language of the other child of the concatenation contains the empty string, which is already known as a bit-flag during construction. For example, the regex  $\mathbf{b} \backslash \mathbf{A}$  is trivially unsatisfiable, and rewritten to  $\perp$ , as the language of  $\mathbf{b}$  does not contain the empty string, therefore the anchor  $\backslash \mathbf{A}$  cannot be at the beginning. Similarly, the regex  $\backslash \mathbf{z} \mathbf{b}$  is rewritten to  $\perp$ , as the language of  $\mathbf{b}$  does not contain the empty string, therefore the anchor  $\backslash \mathbf{z}$  cannot be at the end.

The second rule in Figure 5 provides a more efficient way to check for satisfiability of the complement operation, as the only unsatisfiable complement is  $\sim(\_*)$ . This means that we can check that a complement is satisfiable by ensuring that the complement body is not equivalent to  $\_*$ , effectively searching for the existence of a non-accepting state in the body of the complement. While it is not more efficient asymptotically, it is slightly more efficient memory-wise, as we do not need to construct derivatives of the (outer) complement node itself. Often,

in practice, such non-accepting state is the very state we start from, skipping the check entirely.

In the solver implementation, we do not perform the satisfiability check for the intersection operation eagerly, as the tasks often consist of many intersections, and proving nonemptiness of a subset of intersections does not necessarily help us in proving nonemptiness of all the intersections combined. Instead, we label each regex with a bit-flag, that indicates whether the regex contains an intersection, and only perform the satisfiability check when requested. If the flag is missing, we can be sure that the regex is already known to be satisfiable.

The solver also includes several variants of rewrite rules that are applied to *bounded loops*. In particular, it uses the LOOP rule and the predicate subsumption rule from [46, Figure 6]. To summarize, the majority of time spent in the solver is proving nonemptiness of intersection, which is done lazily, while the rest of the operations are performed eagerly during construction.

## 6 Evaluation

The evaluation is based on existing SMT benchmarks. All regular expressions supported in SMT belong to the class **ERE** and therefore do not support anchors, i.e., in terms of **ERE<sub>↓</sub>**,  $R$  can be treated as  $\backslash AR \backslash z$ . For language semantics in **ERE<sub>↓</sub>**, anchors can also be systematically eliminated through preprocessing. We did not use such preprocessing here, because our evaluation did not require it.

We have collected a set of SMT benchmarks for regexes. The benchmarks consist of approximately 20 000 SMT-LIB files, which are a mix of satisfiability problems for regexes, such as equivalence, subsumption, and membership queries. The benchmarks include AutomataArk [21], Denghang and Sygus-qgen benchmark sets from the QF\_S and QF\_SLIA categories of the SMT-LIB [6] repository, as well as additional benchmarks from the repository [42] designed specifically around difficult regular expression problems.

We compared our implementation (**ERE#-solver**) with the following state-of-the-art tools: cvc5 [5], Z3 [49], Z3str3 [10], Z3str4 [36], Z3alpha [34], OSTRICH [17], OSTRICH<sup>RECL</sup> [27,28] and Z3-Noodler [19]. The benchmark results (when solved) are identical between **ERE#-solver**, cvc5, Z3 and Z3-Noodler. There are discrepancies in the results of some solvers that we have not yet investigated. Table 1 shows the number of unsolved problems for each benchmark category.

Our solver is a specialized tool for regex constraint solving and analysis, rather than a general purpose SMT solver. As a result, it only supports a subset of the SMT operations that the other tools support. Nevertheless, for the regex benchmarks, our solver consistently outperforms the other tools, see Figure 1 in the introduction for the cactus plot of the time taken to solve the benchmarks, and Table 1 for the number of unsolved benchmarks.

For our solver, under 100 of the benchmarks took more than 0.01s to solve, and none of the benchmarks took more than 0.1s to solve. Many individual benchmarks take very little time to solve for other solvers as well, e.g., cvc5

Table 1: Unsolved problems by solver in the various benchmark categories, where **min** and **max** are highlighted in color. Timeout for each problem is **6s**.

	Syg	Denghang	AutArk	Bool	Date	Blowup	Passwd	Mem	Int	Sub	State	$\Sigma$
<b>Included</b>	<b>343</b>	<b>999</b>	<b>15995</b>	<b>21</b>	<b>19</b>	<b>14</b>	<b>34</b>	<b>1907</b>	<b>55</b>	<b>100</b>	<b>22</b>	<b>19509</b>
ERE#-solver	0	0	0	0	0	0	0	0	0	0	0	0
Z3-Noodler	0	1	41	0	0	0	0	0	0	0	0	42
OSTRICH	0	0	75	0	0	0	1	1	2	3	0	82
OSTR <sup>RECL</sup>	0	26	76	0	0	0	0	1	1	3	1	108
cvc5	1	26	152	2	3	1	7	3	27	28	0	250
Z3	0	124	322	1	0	2	16	2	20	9	8	504
Z3alpha	0	126	210	1	0	3	18	1537	55	99	8	2057
Z3str4	0	3	58	1	5	1	17	1907	55	100	6	2153
Z3str3	<b>37</b>	<b>773</b>	<b>5353</b>	1	<b>14</b>	0	<b>20</b>	129	21	26	7	<b>6381</b>

solves thousands of the benchmarks in less than 0.01s, and is generally closest to the time of our solver in Figure 1. The high lower-bound of the time taken for OSTRICH is likely due to a cold start of the process, as it is a Java-based tool, but in terms of the number of benchmarks solved, it is the third best tool after our solver and Z3-Noodler. The key factors contributing to our solver’s significant performance improvement are discussed in the following section.

**Performance Analysis.** The key differences between our solver and others are in the way in which we represent and solve the problem:

- Our solver contains the problem fully within the regex theory, and does not need to convert any part of the problem to a different theory, such as linear arithmetic. Several other solvers solve regex problems with a combination of regex and linear arithmetic, e.g., encoding length constraints [19].
- Our solver never constructs an automaton for the regex, neither as a DFA nor as an NFA. Instead, it performs the operations symbolically, and the nonemptiness check is done through lazy unfolding of the problem encoded in our regex theory. This is in contrast to many other solvers, which construct an automaton for the regex and perform operations on the automaton.
- Our solver explores *language complement* lazily through derivatives, which is not done in other derivative-based string solvers such as cvc5. Complement is often the bottleneck in equivalence and subsumption checks.

The time required for checking nonemptiness of a regex is *proportional to the number of derivatives constructed*, where each derivative is a deterministic memoized operation that is cached upon the first computation. This reduces the amount of work done to the number of *unique* derivatives reached.

SMT benchmarks for regexes typically fall into two categories: they are either solved within a few milliseconds or remain unsolved due to excessive memory consumption, often caused by an exponential blowup in automaton size. However, our tool avoids the explicit automaton construction, which mitigates the memory usage and allows us to solve many of the benchmarks that other tools fail to solve.

Table 2: Effect of lazy depth-first search on satisfiable problems

Satisfiability Problem	DFS derivatives	BFS derivatives
$(\_a\_)\{25\}&(\_b\_)\{25\}&\{0,50\}$	120	6749
$(\_a\_)\{50\}&(\_b\_)\{50\}&\{0,100\}$	245	47874
$(\_a\_)\{100\}&(\_b\_)\{100\}&\{0,200\}$	495	358249

It is important to note that even the lazy graphs constructed by unwinding the problem can grow exponentially in size, and we have to perform simplifications whenever possible to reduce the effect of such blowup. It is done by rewriting the regexes, both algebraically and on the language level, as described in Section 4.1.

Many benchmarks are specifically designed to expose the exponential blowup. However, they can be solved almost immediately by performing certain simplifications and short-circuiting. In the following we explain the simplifications and their effect on the benchmarks in more detail.

**State Space Considerations.** Although the graphs can grow exponentially in size, our solver has several optimizations to mitigate blowup. In many practical cases where other solvers run out of memory, our heuristics enable us to solve the problems much more efficiently. Here we highlight some key optimizations and heuristics demonstrated on problems from the set of collected benchmarks.

Many benchmarks where solvers run out of memory involve the use of large quantifiers, such as `str.len x > 1000` or `re.loop 100 100`. These come with an enormous growth of state space, but as we *do not construct an automaton*, and are not looking for the shortest solution, a depth-first traversal of the state graph is sufficient. Table 2 shows the impact of the approach on satisfiable benchmarks with large quantifiers. The benchmarks are measured by the number of derivatives constructed to reach a solution, with depth-first search significantly outperforming breadth-first search. E.g., depth-first search solves the last benchmark in Table 2 with 495 derivatives in less than 0.002s, including parsing and pre-processing the problem itself, while breadth-first search takes about 0.725s.

While it has an impressive effect on satisfiable problems with a large number of possible solutions, it brings limited advantage for unsatisfiable benchmarks, as search has to exhaust all possibilities before concluding that the regex is unsatisfiable. For unsatisfiable problems, early elimination of trivially unsatisfiable regexes is crucial, as it can reduce search time by several orders of magnitude.

A key simplification for unsatisfiable problems is to prioritize checking length constraints of intersections, as this can be done without taking derivatives and will often eliminate the need for further search.

For example, the regex  $\_ \{4000,5000\} \& \_ \{8000,9000\}$  can be rewritten to  $\perp$  in constant time, as there is no overlap between the minimum and maximum lengths of the two. If there is an overlap, the regex can be rewritten to the intersection of the two length constraints, as was shown in Section 5.

Table 3: Effect of solving unsat suffix problems in reverse

Satisfiability Problem	Rev. derivatives	Fwd. derivatives
<code>*b_{10}&amp;*_a_{10}&amp;_{10,}abc_{10,}</code>	9	279
<code>*b_{20}&amp;*_a_{20}&amp;_{20,}abc_{20,}</code>	19	11 833
<code>*b_{30}&amp;*_a_{30}&amp;_{30,}abc_{30,}</code>	29	539 145
<code>*b_{40}&amp;*_a_{40}&amp;_{40,}abc_{40,}</code>	39	5 000 000+

Another important heuristic is checking the emptiness of a regex in reverse. It is particularly useful for addressing a common pathological case where large state space is combined with an unsatisfiable suffix. For example, the satisfiability of the regex `*b_{100}&*_a_{100}` can be checked for emptiness by reversing the regex to `_{100}b_*&_{100}a_*` and checking it for emptiness. The reversed approach simplifies the problem significantly, as illustrated in Table 3.

## 7 Proposal for SMT-LIB

We believe that a conservative extension of SMT-LIB to support  $\text{RegLan} = \mathbf{ERE}_{\leq}$  would be noncontroversial. For positive and negative lookaheads, positive and negative lookbehinds, and *span membership* or *matching* the minimal extension could be as follows where the start anchor `re.A` and the end anchor `re.z` could be included as regex constants for convenience:

$$\begin{aligned}
(\text{re.lb } r) &\stackrel{\text{DEF}}{=} (?<=r) & (\text{re.nlb } r) &\stackrel{\text{DEF}}{=} (?<!r) & \text{re.A} &\stackrel{\text{DEF}}{=} (\text{re.nlb re.allchar}) \\
(\text{re.la } r) &\stackrel{\text{DEF}}{=} (?=r) & (\text{re.nla } r) &\stackrel{\text{DEF}}{=} (?!r) & \text{re.z} &\stackrel{\text{DEF}}{=} (\text{re.nla re.allchar}) \\
(\text{str.matches } u \ v \ s \ r) &\stackrel{\text{DEF}}{=} (u, v, s) \models r
\end{aligned}$$

Regarding compatibility with  $\mathbf{ERE}$ , `(str.in_re v r)` fully retains its original interpretation for  $r \in \mathbf{ERE}$ . It follows by easy induction over  $r \in \mathbf{ERE}$  that  $\forall u, v, s : (u, v, s) \models r \Leftrightarrow v \in \mathcal{L}(r)$ , i.e.,  $v \in \text{match}(\mathcal{M}(r)) \Leftrightarrow v \in \mathcal{L}(r)$ . Thus,

$$\forall u, v, s \in \text{String}, r \in \mathbf{ERE} : (\text{str.matches } u \ v \ s \ r) \Leftrightarrow (\text{str.in\_re } v \ r)$$

Using lookarounds, one can support other standard anchors such as the *line* anchors `^` and `$` and the *wordborder* anchor `\b` as `(?<!\w)(?=\w) | (?<=\w)(?! \w)`. One can also express different match semantics like *POSIX match* semantics of  $w = uvs \wedge (u, v, s) \models r$ , where  $|u|$  is minimal and  $|v|$  maximal for  $|u|$  – the *leftmost longest* match in  $w$  – as an SMT optimization problem.

It is also practical to support  $\oplus$ , say `re.xor`, and  $\odot$ , say `re.xnor`, natively, in order to avoid their indirect encodings in various algorithms, such as equivalence checking, since  $r_1 \equiv r_2 \Leftrightarrow r_1 \oplus r_2 \equiv \perp$ . Moreover, *reversal* of sequences and regexes is beneficial, as illustrated above. In particular,  $(u, v, s) \models r \Leftrightarrow (s^r, v^r, u^r) \models r^r$ .

**Acknowledgement.** E. Zhuchko is supported by the Estonian Research Council grant *Automata in Learning, Interaction and Concurrency* (PRG1210).

**Disclosure of Interests.** The authors have no competing interests to declare that are relevant to the content of this article.

## References

1. Almeida, M., Moreira, N., Reis, R.: Antimirov and Mosses’s rewrite system revisited. *Int. J. Found. Comput. Sci.* **20**(4), 669–684 (2009), <https://doi.org/10.1142/S0129054109006802>
2. Almeida, M., Moreira, N., Reis, R.: Testing the equivalence of regular languages. *J. Autom. Lang. Comb.* **15**(1/2), 7–25 (2010), <https://doi.org/10.25596/jalc-2010-007>
3. Amadini, R.: A survey on string constraint solving. *ACM Comput. Surv.* **55**(2), 16:1–16:38 (2023), <https://doi.org/10.1145/3484198>
4. Antimirov, V.M.: Partial derivatives of regular expressions and finite automaton constructions. *Theor. Comput. Sci.* **155**(2), 291–319 (1996), [https://doi.org/10.1016/0304-3975\(95\)00182-4](https://doi.org/10.1016/0304-3975(95)00182-4)
5. Barbosa, H., Barrett, C.W., Brain, M., Kremer, G., Lachnitt, H., Mann, M., Mohamed, A., Mohamed, M., Niemetz, A., Nötzli, A., Ozdemir, A., Preiner, M., Reynolds, A., Sheng, Y., Tinelli, C., Zohar, Y.: cvc5: A versatile and industrial-strength SMT solver. In: Fisman, D., Rosu, G. (eds.) *TACAS. LNCS*, vol. 13243, pp. 415–442. Springer (2022). [https://doi.org/10.1007/978-3-030-99524-9\\_24](https://doi.org/10.1007/978-3-030-99524-9_24)
6. Barrett, C., Fontaine, P., Tinelli, C.: The Satisfiability Modulo Theories Library (SMT-LIB) (2016), <https://www.SMT-LIB.org>
7. Barrière, A., Pit-Claudel, C.: Linear matching of javascript regular expressions. *Proc. ACM Program. Lang.* **8** (Jun 2024). <https://doi.org/10.1145/3656431>
8. Berglund, M., van der Merwe, B., van Litsenborgh, S.: Regular expressions with lookahead. *Journal of Universal Computer Science* **27**(4), 324–340 (2021). <https://doi.org/10.3897/jucs.66330>
9. Berzish, M., Day, J.D., Ganesh, V., Kulczynski, M., Manea, F., Mora, F., Nowotka, D.: Towards more efficient methods for solving regular-expression heavy string constraints. *Theor. Comput. Sci.* **943**, 50–72 (2023), <https://doi.org/10.1016/j.tcs.2022.12.009>
10. Berzish, M., Ganesh, V., Zheng, Y.: Z3str3: A string solver with theory-aware heuristics. In: Stewart, D., Weissenbacher, G. (eds.) *FMCAD*. pp. 55–59. IEEE (2017). <https://doi.org/10.23919/FMCAD.2017.8102241>
11. Berzish, M., Kulczynski, M., Mora, F., Manea, F., Day, J.D., Nowotka, D., Ganesh, V.: An SMT solver for regular expressions and linear arithmetic over string length. In: Silva, A., Leino, K.R.M. (eds.) *CAV. LNCS*, vol. 12760, pp. 289–312. Springer (2021). [https://doi.org/10.1007/978-3-030-81688-9\\_14](https://doi.org/10.1007/978-3-030-81688-9_14)
12. Bjørner, N., Fazekas, K.: On incremental pre-processing for smt. In: Pientka, B., Tinelli, C. (eds.) *Automated Deduction – CADE 29*. pp. 41–60. Springer, Cham (2023). [https://doi.org/10.1007/978-3-031-38499-8\\_3](https://doi.org/10.1007/978-3-031-38499-8_3)
13. Bjørner, N., Ganesh, V., Michel, R., Veanes, M.: An SMT-LIB format for sequences and regular expressions. In: Fontaine, P., Goel, A. (eds.) *SMT*. pp. 76–86 (2012). <https://doi.org/10.29007/w5m5>
14. Blahoudek, F., Chen, Y., Chocholatý, D., Havlena, V., Holík, L., Lengál, O., Síc, J.: Word equations in synergy with regular constraints. In: Chechik, M., Katoen, J., Leucker, M. (eds.) *Formal Methods. LNCS*, vol. 14000, pp. 403–423. Springer (2023). [https://doi.org/10.1007/978-3-031-27481-7\\_23](https://doi.org/10.1007/978-3-031-27481-7_23)
15. Brzozowski, J.A.: Derivatives of regular expressions. *JACM* **11**, 481–494 (1964). <https://doi.org/10.1145/321239.321249>
16. Chen, T., Flores-Lamas, A., Hague, M., Han, Z., Hu, D., Kan, S., Lin, A.W., Rümmer, P., Wu, Z.: Solving string constraints with regex-dependent functions

- through transducers with priorities and variables. *Proc. ACM Program. Lang.* **6**(POPL) (jan 2022), <https://doi.org/10.1145/3498707>
17. Chen, T., Hague, M., Lin, A.W., Rümmer, P., Wu, Z.: Decision procedures for path feasibility of string-manipulating programs with complex operations. *Proc. ACM Program. Lang.* **3**(POPL), 49:1–49:30 (2019), <https://doi.org/10.1145/3290362>
  18. Chen, Y., Chocholatý, D., Havlena, V., Holík, L., Lengál, O., Síc, J.: Solving string constraints with lengths by stabilization. *Proc. ACM Program. Lang.* **7**(OOPSLA2), 2112–2141 (2023), <https://doi.org/10.1145/3622872>
  19. Chen, Y., Chocholatý, D., Havlena, V., Holík, L., Lengál, O., Síc, J.: Z3-noodler: An automata-based string solver. In: Finkbeiner, B., Kovács, L. (eds.) TACAS. LNCS, vol. 14570, pp. 24–33. Springer (2024). [https://doi.org/10.1007/978-3-031-57246-3\\_2](https://doi.org/10.1007/978-3-031-57246-3_2)
  20. Chocholatý, D., Fiedor, T., Havlena, V., Holík, L., Hruska, M., Lengál, O., Síc, J.: Mata: A fast and simple finite automata library. In: Finkbeiner, B., Kovács, L. (eds.) TACAS. LNCS, vol. 14571, pp. 130–151. Springer (2024). [https://doi.org/10.1007/978-3-031-57249-4\\_7](https://doi.org/10.1007/978-3-031-57249-4_7)
  21. D’Antoni, L.: AutomataArk (2018), <https://github.com/lorisdanto/automatark>
  22. D’Antoni, L., Veanes, M.: Automata modulo theories. *Commun. ACM* **64**(5), 86–95 (2021), <https://doi.org/10.1145/3419404>
  23. Eriksson, B., Stjerna, A., Masellis, R.D., Rümmer, P., Sabelfeld, A.: Black ostrich: Web application scanning with string solvers. In: Meng, W., Jensen, C.D., Cremers, C., Kirda, E. (eds.) CCS. pp. 549–563. ACM (2023). <https://doi.org/10.1145/3576915.3616582>
  24. Gallant, A.: Burntsushi: rebar (2024), <https://github.com/BurntSushi/rebar>
  25. Hague, M., Jez, A., Lin, A.W.: Parikh’s theorem made symbolic. *Proc. ACM Program. Lang.* **8**(POPL), 1945–1977 (2024), <https://doi.org/10.1145/3632907>
  26. Hojjat, H., Rümmer, P., Shamakhi, A.: On strings in software model checking. In: Lin, A.W. (ed.) APLAS. LNCS, vol. 11893, pp. 19–30. Springer (2019). [https://doi.org/10.1007/978-3-030-34175-6\\_2](https://doi.org/10.1007/978-3-030-34175-6_2)
  27. Hu, D., Wu, Z.: String constraints with regex-counting and string-length solved more efficiently. In: Hermanns, H., Sun, J., Bu, L. (eds.) SETTA. LNCS, vol. 14464, pp. 1–20. Springer (2023), [https://doi.org/10.1007/978-981-99-8664-4\\_1](https://doi.org/10.1007/978-981-99-8664-4_1)
  28. Hu, D., Wu, Z.: An efficient string solver for string constraints with regex-counting and string-length. *Journal of Systems Architecture* pp. 1–38 (2025), <https://doi.org/10.1016/j.sysarc.2025.103340>
  29. Jez, A., Lin, A.W., Markgraf, O., Rümmer, P.: Decision procedures for sequence theories. In: Enea, C., Lal, A. (eds.) CAV. LNCS, vol. 13965, pp. 18–40. Springer (2023). [https://doi.org/10.1007/978-3-031-37703-7\\_2](https://doi.org/10.1007/978-3-031-37703-7_2)
  30. Le, Q.L., He, M.: A decision procedure for string logic with quadratic equations, regular expressions and length constraints. In: Ryu, S. (ed.) APLAS. LNCS, vol. 11275, pp. 350–372. Springer (2018). [https://doi.org/10.1007/978-3-030-02768-1\\_19](https://doi.org/10.1007/978-3-030-02768-1_19)
  31. Liang, T., Reynolds, A., Tsiskaridze, N., Tinelli, C., Barrett, C.W., Deters, M.: An efficient SMT solver for string constraints. *Formal Methods Syst. Des.* **48**(3), 206–234 (2016), <https://doi.org/10.1007/s10703-016-0247-6>
  32. Liang, T., Tsiskaridze, N., Reynolds, A., Tinelli, C., Barrett, C.W.: A decision procedure for regular membership and length constraints over unbounded strings. In: Lutz, C., Ranise, S. (eds.) FroCoS. LNCS, vol. 9322, pp. 135–150. Springer (2015). [https://doi.org/10.1007/978-3-319-24246-0\\_9](https://doi.org/10.1007/978-3-319-24246-0_9)

33. Lotz, K., Goel, A., Dutertre, B., Kiesl-Reiter, B., Kong, S., Majumdar, R., Nowotka, D.: Solving string constraints using SAT. In: Enea, C., Lal, A. (eds.) CAV. LNCS, vol. 13965, pp. 187–208. Springer (2023). [https://doi.org/10.1007/978-3-031-37703-7\\_9](https://doi.org/10.1007/978-3-031-37703-7_9)
34. Lu, Z., Siemer, S., Jha, P., Day, J.D., Manea, F., Ganesh, V.: Layered and staged monte carlo tree search for SMT strategy synthesis. In: IJCAI. pp. 1907–1915. ijcai.org (2024). <https://doi.org/10.24963/ijcai.2024/211>
35. Mamouras, K., Chattopadhyay, A.: Efficient matching of regular expressions with lookaround assertions. Proc. ACM Program. Lang. **8**(POPL), 2761–2791 (2024), <https://doi.org/10.1145/3632934>
36. Mora, F., Berzish, M., Kulczynski, M., Nowotka, D., Ganesh, V.: Z3str4: A multi-armed string solver. In: FM. pp. 389–406. Springer, Berlin, Heidelberg (2021). [https://doi.org/10.1007/978-3-030-90870-6\\_21](https://doi.org/10.1007/978-3-030-90870-6_21)
37. Moreira, N., Pereira, D., de Sousa, S.M.: Deciding regular expressions (in-)equivalence in coq. In: Kahl, W., Griffin, T.G. (eds.) RAMiCS. LNCS, vol. 7560, pp. 98–113. Springer, Berlin, Heidelberg (2012). [https://doi.org/10.1007/978-3-642-33314-9\\_7](https://doi.org/10.1007/978-3-642-33314-9_7)
38. Moseley, D., Nishio, M., Perez Rodriguez, J., Saarikivi, O., Toub, S., Veanes, M., Wan, T., Xu, E.: Derivative based nonbacktracking real-world regex matching with backtracking semantics. Proc. ACM Program. Lang. **7**(PLDI), 148:1–148:24 (Jun 2023). <https://doi.org/10.1145/3591262>
39. de Moura, L., Bjørner, N.: Z3: An Efficient SMT Solver. In: TACAS. LNCS, vol. 4963, pp. 337–340. Springer, Berlin, Heidelberg (2008). [https://doi.org/10.1007/978-3-540-78800-3\\_24](https://doi.org/10.1007/978-3-540-78800-3_24)
40. Owens, S., Reppy, J.H., Turon, A.: Regular-expression derivatives re-examined. J. Funct. Program. **19**(2), 173–190 (2009), <https://doi.org/10.1017/S0956796808007090>
41. Reynolds, A., Woo, M., Barrett, C.W., Brumley, D., Liang, T., Tinelli, C.: Scaling up DPLL(T) string solvers using context-dependent simplification. In: Majumdar, R., Kuncak, V. (eds.) CAV. LNCS, vol. 10427, pp. 453–474. Springer (2017). [https://doi.org/10.1007/978-3-319-63390-9\\_24](https://doi.org/10.1007/978-3-319-63390-9_24)
42. Stanford, C.: Boolean Regular Expression SMT Benchmarks (2021), <https://github.com/cdstanford/regex-smt-benchmarks>
43. Stanford, C., Veanes, M., Bjørner, N.S.: Symbolic boolean derivatives for efficiently solving extended regular expression constraints. In: Freund, S.N., Yahav, E. (eds.) PLDI. pp. 620–635. ACM (2021). <https://doi.org/10.1145/3453483.3454066>
44. Turoňová, L., Holík, L., Homoliak, I., Lengál, O., Veanes, M., Vojnar, T.: Counting in regexes considered harmful: Exposing ReDoS vulnerability of nonbacktracking matchers. In: 31st USENIX Security Symposium (USENIX Security 22). pp. 4165–4182. USENIX Association, Boston, MA (August 2022), <https://www.usenix.org/conference/usenixsecurity22/presentation/turonova>
45. Varatalu, I.E.: Implementation of RE#-solver (2025), <https://github.com/ieviev/cav25-resharp-smt>
46. Varatalu, I.E., Veanes, M., Ernits, J.: RE#: High performance derivative-based regex matching with intersection, complement, and restricted lookarounds. Proc. ACM Program. Lang. **9**(POPL), 1:1–1:32 (January 2025), <https://doi.org/10.1145/3704837>
47. Varatalu, I.E., Zhuchko, E., Ernits, J.: Artifact for regex decision procedures in extended RE# (2025), <https://doi.org/10.5281/zenodo.15210805>

48. Veanes, M., Ball, T., Ebner, G., Zhuchko, E.: Symbolic automata: Omega-regularity modulo theories. *Proc. ACM Program. Lang.* **9**(POPL), 2:1–2:34 (January 2025), <https://doi.org/10.1145/3704838>
49. Z3 Prover <https://github.com/Z3Prover/z3/wiki>
50. Zhuchko, E., Veanes, M., Ebner, G.: Lean formalization of extended regular expression matching with lookarounds. In: Timany, A., Traytel, D., Pientka, B., Blazy, S. (eds.) *CPP*. pp. 118–131. ACM (2024). <https://doi.org/10.1145/3636501.3636959>