

Verified decision procedures for MSO on words based on derivatives of regular expressions

DMITRIY TRAYTEL

Institute of Information Security, Department of Computer Science, ETH Zürich, Switzerland
(e-mail: traytel@inf.ethz.ch)

TOBIAS NIPKOW

Fakultät für Informatik, Technische Universität München, Germany
(e-mail: nipkow@in.tum.de)

Abstract

Monadic second-order logic on finite words is a decidable yet expressive logic into which many decision problems can be encoded. Since MSO formulas correspond to regular languages, equivalence of MSO formulas can be reduced to the equivalence of some regular structures (e.g., automata). This paper presents a verified functional decision procedure for MSO formulas that is not based on automata but on regular expressions. Functional languages are ideally suited for this task: regular expressions are data types and functions on them are defined by pattern matching and recursion and are verified by structural induction. Decision procedures for regular expression equivalence have been formalized before, usually based on Brzozowski derivatives. Yet, for a straightforward embedding of MSO formulas into regular expressions, an extension of regular expressions with a projection operation is required. We prove total correctness and completeness of an equivalence checker for regular expressions extended in that way. We also define a language-preserving translation of formulas into regular expressions with respect to two different semantics of MSO. Our results have been formalized and verified in the theorem prover Isabelle. Using Isabelle's code generation facility, this yields purely functional, formally verified programs that decide equivalence of MSO formulas.

1 Introduction

Many decision procedures for logical theories are based on the famous logic-automaton connection. That is, they reduce the decision problem for some logical theory to a decidable question about some class of automata. Automata are usually implemented with the help of imperative data structures for efficiency reasons.

In functional languages, automata are not an ideal abstraction because they are graphs rather than trees. In contrast, regular expressions are perfect for functional languages and they are equally expressive. In fact, Brzozowski (1964) showed how automata-based algorithms can be recast as recursive algebraic manipulations of regular expressions. His *derivatives* can be seen as a way of simulating automaton states with regular expressions and computing the next-state function symbolically.

Recently, Brzozowski's derivatives were discovered by functional programmers and theorem provers. Owens *et al.* (2009) realized that regular expressions and

their derivatives fit perfectly with data types and recursive functions. Their paper explores regular expression matching based directly on regular expressions rather than automata. Fischer *et al.* (2010) also explore regular expression matching, but by means of marked regular expressions rather than derivatives. Slightly later, the interactive theorem proving community woke up to the beauty of derivatives, too. This resulted in four papers about verified decision procedures for the equivalence of regular expressions based on derivatives and on marked regular expressions (see related work below). In one of these four papers, Coquand and Siles (2011) state that “A more ambitious project will be to use this work for writing a decision procedure for WS1S”, a monadic second-order (MSO) logic. Our paper does just that (and more).

MSO logic on finite words is a decidable yet expressive logic into which many decision problems can be encoded (Thomas, 1997). MSO allows only monadic predicates but quantification both over numbers and finite sets of numbers. Two closely related but subtly different semantics can be found in the literature. One of the two, WS1S—the Weak MSO logic of 1 Successor, is based on arithmetic. The other, M2L(Str) (Henriksen *et al.*, 1995), is more closely related to formal languages. There seems to be some disagreement as to which semantics is the more appropriate one for verification purposes (Klarlund, 1999; Ayari & Basin, 2000). Hence, we cover both.

Essentially, MSO formulas describe regular languages. Therefore, MSO formulas can be decided by translating them into automata. This is the basis of the highly successful MONA tool (Elgaard *et al.*, 1998) for deciding WS1S. MONA’s success is due to its (in practical terms) highly efficient implementation and to the ease with which very different verification problems can be encoded in MSO logic, for example Presburger arithmetic and Hoare logic for pointer programs.

The contribution of this paper is the presentation of the first purely functional decision procedures for two interpretations of MSO based on derivatives of regular expressions. These decision procedures have been verified in Isabelle/HOL and we sketch their correctness proofs. We are not aware of any previous decision procedure for MSO based on regular expressions (as opposed to automata), let alone a verified program.

It is instructive to compare our decision procedure for WS1S with MONA. MONA is a highly tuned implementation using cache-conscious data structures including a BDD-based automaton representation. Ours is a (by comparison tiny) purely functional program that operates on regular expressions and can only cope with small examples. MONA is not verified (and the prospect of doing so is daunting), whereas our code is.

In this paper, we distinguish ordinary regular expressions that contain only concatenation, union, and iteration from *extended* regular expressions that also provide complement and intersection. The rest of the paper is organized as follows. Section 2 gives an overview of related work. Section 3 introduces some basic notations. Sections 4 and 5 constitute the main contribution of our paper—the first shows how to decide equivalence of extended regular expressions with an additional projection operation, the second reduces equivalence of MSO formulas to equivalence of exactly those regular expressions with respect to both semantics,

M2L and WS1S. In total, this yields a decision procedure for MSO on words. A short case study of the decision procedure is given in Section 6.

This paper is an extended and revised version of the homonymous ICFP 2013 functional pearl (Traytel & Nipkow, 2013). The new contributions are mostly actual changes to the decision procedure aiming to improve on both, performance and presentation:

- Picking up an idea from Owens *et al.* (2009), we change the semantics of atomic regular expressions to represent sets of alphabet letters rather than single letters (Section 4.4).
- Based on earlier experimental results (Nipkow & Traytel, 2014), we change the backend decision procedure for regular expression equivalence to use partial derivatives instead of Brzowski derivatives (Section 4.5).
- We improve the translations of M2L(Str) formulas to regular expressions by removing redundancies and expand on the previously omitted implementation of the translation of WS1S formulas to regular expressions (Section 5.4).
- A reevaluation of the performance shows a sizable improvement, yet still far away from competing with MONA (Section 6).

While the paper is intended to be self-contained with respect to the presented functional program deciding equivalence of MSO formulas, we deliberately give only rough intuitions instead of detailed proofs. The proofs are where they truly belong: in the publicly available formalization (Traytel & Nipkow, 2014).

2 Related work

Brzowski (1964) introduced the notion of derivatives of extended regular expressions and Ginzburg (1967) employed them in an algorithm for deciding language equivalence that we essentially are using here. Antimirov (1996) devised the related notion of partial derivatives of ordinary regular expressions. Caron *et al.* (2011) extended partial derivatives to extended regular expressions. The concept of derivatives as means to compute the next state symbolically goes beyond regular expressions— as witnessed by libraries for parsing developed by Danielsson (2010) in Agda and by Might *et al.* (2011) in Lisp using lazily evaluated variations of Brzowski derivatives for parser combinators. Furthermore, Kozen (2008) lifted derivatives to expressions of Kleene algebra with tests.

MONA was linked to Isabelle by Basin & Friedrich (2000) and to PVS by Owre and Rueß (2000). In both cases, MONA is used as a trusted oracle for deciding formulas in the respective theorem prover.

Now, we discuss work on verified decision procedures for regular expressions. The first verified equivalence checker for regular expressions was published by Braibant & Pous (2010). They worked with automata, not regular expressions, their theory was large and their algorithm efficient. In response, Krauss & Nipkow (2012) gave a much simpler partial correctness proof for an equivalence checker for regular expressions based on derivatives. Coquand & Siles (2011) showed total correctness of their equivalence checker for extended regular expressions based on

derivatives. Asperti (2012) presented an equivalence checker for regular expressions via marked regular expressions (as previously used by Fischer *et al.* (2010)) and showed total correctness. Moreira *et al.* (2012) presented an equivalence checker for regular expressions based on partial derivatives and showed its total correctness. Recently, we have devised a general framework that unifies the different approaches based on derivatives, partial derivatives, and marked regular expressions under one roof (Nipkow & Traytel, 2014). Berghofer & Reiter (2009) formalized a decision procedure for Presburger arithmetic via automata in Isabelle/HOL.

Outside of the application area of equivalence checking, Wu *et al.* (2014) benefited from the inductive structure of regular expressions to formally verify the Myhill–Nerode theorem.

3 Preliminaries

Although we formalized everything in this paper in the theorem prover Isabelle/HOL (Nipkow *et al.*, 2002; Nipkow & Klein, 2014), no knowledge of theorem provers or Isabelle/HOL is required because we employ mostly ordinary mathematical notation in our presentation. Some specific notations are summarized below.

The symbol \mathbb{B} represents the type of Booleans, where \top and \perp represent true and false. The type of sets and the type of lists over some type τ are written τ set and τ list. In general, type constructors follow their arguments. The letters α and β represent type variables. The notation $t :: \tau$ means that term t has type τ .

Many of our functions are curried. In some cases, we write the first argument as an index: instead of $f a b$, we write $f_a(b)$ (in preference to just $f_a b$). The projection functions on pairs are called `fst` and `snd`. The image of a function f over a set S is written $f \cdot S$.

Lists are built up from the empty list `[]` via the infix `#` operator that prepends an element x to a list xs : $x \# xs$. Two lists are concatenated with the infix `@` operator. Accessing the n th element of a list xs is denoted by $xs[n]$; the indexing is zero-based. The length of the list xs is written $|xs|$.

Finite words as in formal language theory are modeled as finite lists, i.e., type α list. The empty word is the empty list. As is customary, concatenation of two words u and v is denoted by their juxtaposition uv ; similarly for a single letter a of the alphabet and a word w : aw . That is, the operators `#` and `@` remain implicit (for words, not for arbitrary lists).

4 Extended regular expressions

In Section 5, MSO formulas are translated into regular expressions such that encodings of models of a formula correspond exactly to words in the regular language. Thereby, equivalence of formulas is reduced to the equivalence of regular expressions.

Decision procedures for equivalence of regular expression have been formalized earlier in theorem provers. Here, we extend the existing formalization and the

soundness proof in Isabelle/HOL by Krauss & Nipkow (2012) with negation and intersection operation on regular expressions, as well as with a nonstandard projection operation. Additionally, we provide proofs of termination and completeness.

4.1 Syntax and semantics

Regular expressions extended with intersection and complement allow us to encode Boolean operators on formulas in a straightforward fashion. A further operation—the *projection* Π —plays the crucial role of encoding existential quantifiers. These Π -extended regular expressions (to distinguish them from mere extended regular expressions) are defined as a recursive data type α RE, where α is the type of the underlying alphabet. In conventional concrete syntax, α RE is defined by the grammar

$$\begin{array}{l} r = \mathbf{0} \quad | \quad \mathbf{1} \quad | \quad a \\ \quad | \quad r + s \quad | \quad r \cdot s \quad | \quad r^* \\ \quad | \quad r \cap s \quad | \quad \neg r \quad | \quad \Pi r \end{array}$$

where $r, s :: \alpha$ RE and $a :: \alpha$. Note that much of the time, we will omit the “ Π -extended” and simply speak of regular expressions if there is no danger of confusion.

We assume that type α is partitioned into a family of alphabets Σ_n that depend on a natural number n and there is a function $\pi :: \Sigma_{n+1} \rightarrow \Sigma_n$ ¹ that translates between the different alphabets. In our application, n will represent the number of free variables of the translated MSO formula. For now, Σ_n and π are just parameters of our setup.

We focus on well-formed regular expressions where all atoms come from the same alphabet Σ_n . This will guarantee that the language of such a well-formed expression is a subset of Σ_n^* . The projection operation complicates wellformedness a little. Because projection is meant to encode existential quantifiers, projection should transform a regular expression over Σ_{n+1} into a regular expression over Σ_n , just as the existential quantifier transforms a formula with $n + 1$ free variables into a formula with n free variables. Thus, projection changes the alphabet. Wellformedness is defined as the recursive predicate $\text{wf} :: \mathbb{N} \rightarrow \alpha$ RE $\rightarrow \mathbb{B}$.

$$\begin{array}{ll} \text{wf}_n(\mathbf{0}) & = \top & \text{wf}_n(\mathbf{1}) & = \top \\ \text{wf}_n(a) & = a \in \Sigma_n & \text{wf}_n(r + s) & = \text{wf}_n(r) \wedge \text{wf}_n(s) \\ \text{wf}_n(r \cdot s) & = \text{wf}_n(r) \wedge \text{wf}_n(s) & \text{wf}_n(r^*) & = \text{wf}_n(r) \\ \text{wf}_n(r \cap s) & = \text{wf}_n(r) \wedge \text{wf}_n(s) & \text{wf}_n(\neg r) & = \text{wf}_n(r) \\ \text{wf}_n(\Pi r) & = \text{wf}_{n+1}(r) & & \end{array}$$

We call a regular expression r *n-wellformed* if $\text{wf}_n(r)$ holds.

¹ Due to Isabelle’s lack of dependent types, the actual type of π is $\alpha \rightarrow \alpha$. The more refined dependent type $\Sigma_{n+1} \rightarrow \Sigma_n$ is realized via Isabelle’s tool for modeling parameterized systems with additional assumptions: locales (Ballarín, 2006). A locale fixes parameters and states assumptions about them. Hence, we use the locale assumption $\pi \cdot \Sigma_{n+1} \subseteq \Sigma_n$ to relate locale parameters π and Σ .

The language $\mathcal{L} :: \mathbb{N} \rightarrow \alpha \text{ RE} \rightarrow (\alpha \text{ list}) \text{ set}$ of a regular expression is defined as usual, except for the equations for complement and projection. For an n -well-formed regular expression, the definition yields a subset of Σ_n^* .

$$\begin{aligned} \mathcal{L}_n(\mathbf{0}) &= \{\} & \mathcal{L}_n(\mathbf{1}) &= \{\square\} \\ \mathcal{L}_n(a) &= \{a\} & \mathcal{L}_n(r + s) &= \mathcal{L}_n(r) \cup \mathcal{L}_n(s) \\ \mathcal{L}_n(r \cdot s) &= \mathcal{L}_n(r) \cdot \mathcal{L}_n(s) & \mathcal{L}_n(r^*) &= \mathcal{L}_n(r)^* \\ \mathcal{L}_n(r \cap s) &= \mathcal{L}_n(r) \cap \mathcal{L}_n(s) & \mathcal{L}_n(\neg r) &= \Sigma_n^* \setminus \mathcal{L}_n(r) \\ \mathcal{L}_n(\Pi r) &= \text{map } \pi \cdot \mathcal{L}_{n+1}(r) \end{aligned}$$

The first unusual point is the parametrization with n . It expresses that we expect a regular expression over Σ_n and is necessary for the definition $\mathcal{L}_n(\neg r) = \Sigma_n^* \setminus \mathcal{L}_n(r)$.

The definition $\mathcal{L}_n(\Pi r) = \text{map } \pi \cdot \mathcal{L}_{n+1}(r)$ is parameterized by the fixed parameter $\pi :: \Sigma_{n+1} \rightarrow \Sigma_n$. The projection Π denotes the homomorphic image under π . In more detail: map lifts π homomorphically to words (lists), and \cdot lifts it to sets of words. Therefore, Π transforms a language over Σ_{n+1} into a language over Σ_n .

To understand the “projection” terminology, it is helpful to think of elements of Σ_n as lists of fixed length n over some alphabet Σ and of π as the tail function on lists that drops the first element of the list. A word over Σ_n is then a list of lists. Though this is a good intuition, the actual encoding of formulas later on will be slightly more complicated. Fortunately, we can ignore these complications for now by working with arbitrary but fixed Σ_n and π in the current section. Specific instantiations for them are given in Section 5.

4.2 Deciding language equivalence

Now we turn our attention to deciding equivalence of Π -extended regular expressions. The key concepts required for this are nullability and derivatives. We call a regular expression *nullable* if its language contains the empty word \square . Nullability can be easily checked syntactically by the following recursive function $\varepsilon :: \alpha \text{ RE} \rightarrow \mathbb{B}$.

$$\begin{aligned} \varepsilon(\mathbf{0}) &= \perp & \varepsilon(\mathbf{1}) &= \top \\ \varepsilon(a) &= \perp & \varepsilon(r + s) &= \varepsilon(r) \vee \varepsilon(s) \\ \varepsilon(r \cdot s) &= \varepsilon(r) \wedge \varepsilon(s) & \varepsilon(r^*) &= \top \\ \varepsilon(r \cap s) &= \varepsilon(r) \wedge \varepsilon(s) & \varepsilon(\neg r) &= \neg \varepsilon(r) \\ \varepsilon(\Pi r) &= \varepsilon(r) \end{aligned}$$

The characteristic property— $\varepsilon(r)$ iff $\square \in \mathcal{L}_n(r)$ for any regular expression r and $n :: \mathbb{N}$ —follows by structural induction on r .

The second key concept—the *derivative* of a regular expression $\mathcal{D} :: \alpha \rightarrow \alpha \text{ RE} \rightarrow \alpha \text{ RE}$ and its lifting to words $\mathcal{D}^* :: \alpha \text{ list} \rightarrow \alpha \text{ RE} \rightarrow \alpha \text{ RE}$ —semantically corresponds to left quotients of regular languages with respect to a fixed letter or word. Just as before, the recursive definition is purely syntactic and the semantic correspondence

is established by a straightforward structural induction.

$$\begin{array}{ll}
 \mathcal{D}_b(\mathbf{0}) &= \mathbf{0} & \mathcal{D}_b(\mathbf{1}) &= \mathbf{0} \\
 \mathcal{D}_b(a) &= \text{if } a = b \text{ then } \mathbf{1} \text{ else } \mathbf{0} & \mathcal{D}_b(r + s) &= \mathcal{D}_b(r) + \mathcal{D}_b(s) \\
 \mathcal{D}_b(r \cdot s) &= & \mathcal{D}_b(r^*) &= \mathcal{D}_b(r) \cdot r^* \\
 & \text{if } \varepsilon(r) \text{ then } \mathcal{D}_b(r) \cdot s + \mathcal{D}_b(s) & & \\
 & \text{else } \mathcal{D}_b(r) \cdot s & & \\
 \mathcal{D}_b(r \cap s) &= \mathcal{D}_b(r) \cap \mathcal{D}_b(s) & \mathcal{D}_b(\neg r) &= \neg \mathcal{D}_b(r) \\
 \mathcal{D}_b(\Pi r) &= \Pi \left(\bigoplus_{c \in \pi^{-b}} \mathcal{D}_c(r) \right) & & \\
 \mathcal{D}_{\square}^*(r) &= r & \mathcal{D}_{bw}^*(r) &= \mathcal{D}_w^*(\mathcal{D}_b(r))
 \end{array}$$

Lemma 1

Assume $b \in \Sigma_n$, $v \in \Sigma_n^*$ and let r be an n -well-formed regular expression. Then $\mathcal{L}_n(\mathcal{D}_b(r)) = \{w \mid bw \in \mathcal{L}_n(r)\}$ and $\text{wf}_n(\mathcal{D}_b(r))$, and consequently $\mathcal{L}_n(\mathcal{D}_v^*(r)) = \{w \mid vw \in \mathcal{L}_n(r)\}$ and $\text{wf}_n(\mathcal{D}_v^*(r))$.

The projection case introduced some new syntax that deserves some explanation. The preimage π^- applied to a letter $b \in \Sigma_n$ denotes the set $\{c \in \Sigma_{n+1} \mid \pi c = b\}$. Our alphabets Σ_n are finite for each n , hence so is the preimage of a letter. The summation \bigoplus over a finite set denotes the iterated application of the $+$ -constructor of regular expressions. Summation over the empty set is defined as $\mathbf{0}$.

Derivatives of extended regular expressions were introduced by Brzozowski (1964) fifty years ago. Our contribution is the extension of the concept to handle the projection operation. Since the projection acts homomorphically on words, it is clear that the derivative of Πr with respect to a letter b can be expressed as a projection of derivatives of r . The concrete definition is a consequence of the following identity of left quotients for $b \in \Sigma_n$ and $A \subseteq \Sigma_{n+1}^*$:

$$\{w \mid bw \in \text{map } \pi \cdot A\} = \text{map } \pi \cdot \bigcup_{c \in \pi^{-b}} \{w \mid cw \in A\}$$

Although we completely avoid automata in the formalization, a derivative with respect to the letter b can be seen as a transition labeled by b in a deterministic automaton, the states of which are labeled by regular expressions. The automaton accepting the language of a regular expression r can be thus constructed iteratively by exploring all derivatives of r and defining exactly those states as accepting, which are labeled by a nullable regular expression. However, the set $\{\mathcal{D}_w^*(r) \mid w :: \alpha \text{ list}\}$ of states reachable in this manner is infinite in general. To obtain a finite automaton, the states must be partitioned into classes of regular expressions that are *ACI-equivalent*, i.e., syntactically equal modulo associativity, commutativity and idempotence of the $+$ -constructor or more formally related by the following inductively defined congruence \sim .

$$\begin{array}{lll}
 r + (s + t) \sim (r + s) + t & r + s \sim s + r & r + r \sim r \\
 r \sim r & \frac{r \sim s}{s \sim r} & \frac{r \sim s \quad s \sim t}{r \sim t} \\
 \frac{r_1 \sim s_1 \quad r_2 \sim s_2}{r_1 + r_2 \sim s_1 + s_2} & \frac{r_1 \sim s_1 \quad r_2 \sim s_2}{r_1 \cdot r_2 \sim s_1 \cdot s_2} & \frac{r \sim s}{r^* \sim s^*}
 \end{array}$$

Brzozowski showed that the number of \sim -equivalence classes for a fixed regular expression r is finite by structural induction on r . The inductive steps require proving finiteness by representing equivalence classes of derivatives of the expression in terms of equivalence classes of derivatives of subexpressions. This is technically complicated, especially for concatenation, iteration and projection, since it requires a careful choice of representatives of equivalence classes to reason about them, and Isabelle's automation cannot help much with the finiteness arguments—indeed the verification of Theorem 2 constitutes the most intricate proof in the present work.

Theorem 2

$\{\langle \mathcal{D}_w^*(r) \rangle \mid w :: \alpha \text{ list}\}$ is finite for any regular expression r .

The function $\langle - \rangle :: \alpha \text{ RE} \rightarrow \alpha \text{ RE}$ is the ACI normalization function, which maps ACI-equivalent regular expressions to the same representative, i.e., defines a particular executable choice of representatives of \sim -equivalence classes. It is defined by means of a normalizing constructor $\oplus :: \alpha \text{ RE} \rightarrow \alpha \text{ RE} \rightarrow \alpha \text{ RE}$ and an arbitrary linear order \leq on regular expressions. The equations for \oplus are matched sequentially.

$$\begin{array}{ll}
 \langle \mathbf{0} \rangle & = \mathbf{0} & \langle \mathbf{1} \rangle & = \mathbf{1} \\
 \langle a \rangle & = a & \langle r + s \rangle & = \langle r \rangle \oplus \langle s \rangle \\
 \langle r \cdot s \rangle & = \langle r \rangle \cdot \langle s \rangle & \langle r^* \rangle & = \langle r \rangle^* \\
 \langle r \cap s \rangle & = \langle r \rangle \cap \langle s \rangle & \langle \neg r \rangle & = \neg \langle r \rangle \\
 \langle \Pi r \rangle & = \Pi \langle r \rangle
 \end{array}$$

$$\begin{array}{ll}
 (r + s) \oplus t & = r \oplus (s \oplus t) \\
 r \oplus (s + t) & = \text{if } r = s \text{ then } s + t \\
 & \quad \text{else if } r \leq s \text{ then } r + (s + t) \\
 & \quad \quad \text{else } s + (r \oplus t) \\
 r \oplus s & = \text{if } r = s \text{ then } r \\
 & \quad \text{else if } r \leq s \text{ then } r + s \\
 & \quad \quad \text{else } s + r
 \end{array}$$

When proving Theorem 2 by induction, on a high-level most cases follow Brzozowski's original proof (1964). The only exception is the newly introduced constructor Πr , where we proceed as follows: By induction hypothesis we know that r has a finite set D of distinct derivatives modulo ACI. Some of the formulas in D can have a sum as the topmost constructor. If we repeatedly split such outermost sums in D until none are left, we obtain a finite set X of expressions. Each word derivative $\mathcal{D}_w^*(r)$ is ACI equivalent to some $\Pi (\bigoplus Y)$ for some $Y \subseteq X$. Since X is finite, its powerset is also finite. Hence, there are only finitely many distinct $\mathcal{D}_w^*(r)$ modulo ACI.

The above proof sketch is very informal. The corresponding formal proof is technically more challenging, e.g., we need to define precisely in which way $\mathcal{D}_w^*(r)$ is ACI equivalent to $\Pi (\bigoplus Y)$ for arbitrary words w . Here, we employ the ACI normalization function and its equivalent abstract characterization: After the application of $\langle - \rangle$, all sums in the expression are associated to the right and the summands are

sorted with respect to \leq and duplicated summands are removed. From this, further later on useful properties of $\langle - \rangle$ can be derived:

Lemma 3

Let r be a regular expression, $n :: \mathbb{N}$ and $b \in \Sigma_n$. Then $\mathcal{L}_n \langle r \rangle = \mathcal{L}_n(r)$, $\langle \langle r \rangle \rangle = \langle r \rangle$, and $\langle \mathcal{D}_b \langle r \rangle \rangle = \langle \mathcal{D}_b(r) \rangle$.

So far, ACI normalization only connects Brzozowski derivatives to deterministic finite automata. Furthermore, it will ensure termination of our decision procedure even without ever entering the world of automata. Instead we follow Rutten (1998), who gives an alternative view on deterministic automata as coalgebras. In the coalgebraic setting, the function $\lambda r. (\varepsilon(r), \lambda b. \mathcal{D}_b(r)) :: \alpha \text{ RE} \rightarrow \mathbb{B} \times (\alpha \rightarrow \alpha \text{ RE})$ is a D -coalgebra for the functor $D(S) = \mathbb{B} \times (\alpha \rightarrow S)$. The final coalgebra of D exists and corresponds exactly to the set of all languages. Therefore, we obtain the powerful coinduction principle, reducing language equality to bisimilarity. We phrase this general theorem instantiated to our concrete setting. The formalized proof itself does not require any category theory; it resembles the reasoning in Section 4 of Rutten (1998).

Theorem 4 (Coinduction)

Let $R :: (\alpha \text{ RE} \times \alpha \text{ RE})$ set be a relation, such that for all $(r, s) \in R$, we have the following:

1. $\text{wf}_n(r) \wedge \text{wf}_n(s)$;
2. $\varepsilon(r) \leftrightarrow \varepsilon(s)$;
3. $(\langle \mathcal{D}_b(r) \rangle, \langle \mathcal{D}_b(s) \rangle) \in R$ for all $b \in \Sigma_n$.

Then for all $(r, s) \in R$, $\mathcal{L}_n(r) = \mathcal{L}_n(s)$ holds.

From Lemmas 1 and 3, we know that the relation

$$\mathcal{B} = \{(\langle \mathcal{D}_w^*(r) \rangle, \langle \mathcal{D}_w^*(s) \rangle) \mid w \in \Sigma_n^*\}$$

contains $(\langle r \rangle, \langle s \rangle)$ and fulfills the assumptions 1 and 3 of the coinduction theorem, assuming that r and s are both n -wellformed. Moreover, using Theorem 2, it follows that this relation is finite. Thus, checking assumption 2 for every pair of this finite relation is sufficient to prove language equality of r and s by coinduction. We obtain the following abstract specification of a language equivalence checking algorithm.

Theorem 5

Let r and s be n -well-formed regular expressions. Then $\mathcal{L}_n(r) = \mathcal{L}_n(s)$, iff we have $\varepsilon(r') \leftrightarrow \varepsilon(s')$ for all $(r', s') \in \mathcal{B}$.

4.3 Executable algorithm from a theorem

Our goal is not only to prove some abstract theorems about a decision procedure, but also to extract executable code in some functional programming language (e.g.,

Standard ML, Haskell, OCaml) using the code generation facility of Isabelle/HOL (Haftmann & Nipkow, 2010). Theorem 5 is not enough to do so: it contains a set comprehension ranging over the infinite set Σ_n^* , which is not executable as such. We need to instruct the system how to enumerate \mathcal{B} .

We start with the pair $(\langle r \rangle, \langle s \rangle)$ and compute its pairwise derivatives for all letters of the alphabet. For the computed pairs of regular expressions, we proceed by computing their derivatives and so on. This of course does not terminate. However, if we stop our exploration at pairs that we have seen before it does, since we are exploring a finite set.

In more detail, we use a worklist algorithm that iteratively adds not yet inspected pairs of regular expressions while exhausting words of increasing length until no new pairs are generated. Saturation is reached by means of the executable combinator $\text{while} :: (\alpha \rightarrow \mathbb{B}) \rightarrow (\alpha \rightarrow \alpha) \rightarrow \alpha \rightarrow \alpha$ option from the Isabelle/HOL library. The option type α option has two constructors $\text{None} :: \alpha$ option and $\text{Some} :: \alpha \rightarrow \alpha$ option. Some lifts elements from the base type α to the option type, while None is usually used to indicate some exceptional behavior. The definition of while

$$\text{while } b \ c \ s = \text{if } \exists k. \neg b(c^k(s)) \text{ then } \text{Some } (c^{\text{Least } k. \neg b(c^k(s))}(s)) \text{ else } \text{None}$$

is not executable, but the following key lemma is

$$\text{while } b \ c \ s = \text{if } b \ s \text{ then } \text{while } b \ c \ (c \ s) \text{ else } \text{Some } s$$

The code generated from this recursive equation will return $\text{Some } s$ in case the definition of while says so, but instead of returning None , it will not terminate. Thus, we can prove termination if we can show that the result is $\neq \text{None}$.

In our case, the state s of the while loop consists of a worklist $ws :: (\alpha \text{ RE} \times \alpha \text{ RE})$ list of unprocessed pairs of regular expressions together with a set $N :: (\gamma \times \gamma)$ set of already seen pairs modulo a normalization function $\text{norm} :: \alpha \text{ RE} \rightarrow \gamma$. This normalization function (which is a parameter of our setup) is applied to already ACI-normalized expressions, to syntactically identify further language equivalent expressions. This makes the bisimulation relation that must be exhausted smaller, thus saturation is reached faster. The range type of the normalization is not fixed, but we require a notion of languages $\mathcal{L}^\gamma :: \mathbb{N} \rightarrow \gamma \rightarrow (\alpha \text{ list})$ set to be available for it, such that $\mathcal{L}_n^\gamma(\text{norm } r) = \mathcal{L}_n(r)$ holds. In the simplest case, norm can be the identity function and $\mathcal{L}^\gamma = \mathcal{L}$. More interesting is a function on regular expressions that eliminates $\mathbf{0}$ from unions, concatenations and intersections and $\mathbf{1}$ from concatenations. Other regular structures such as automata or different kinds of regular expressions as instantiations for γ might enable even more sophisticated simplifications.

We define the arguments to the while combinator $b :: (\alpha \text{ RE} \times \alpha \text{ RE})$ list \times $(\gamma \times \gamma)$ set $\rightarrow \mathbb{B}$ and $c :: \mathbb{N} \rightarrow (\alpha \text{ RE} \times \alpha \text{ RE})$ list \times $(\gamma \times \gamma)$ set $\rightarrow (\alpha \text{ RE} \times \alpha \text{ RE})$ list \times $(\gamma \times \gamma)$ set.

$$\begin{aligned} b \ (\[], _) &= \perp \\ b \ ((r, s) \# _, _) &= \varepsilon(r) \leftrightarrow \varepsilon(s) \end{aligned}$$

$$\begin{aligned}
c_n((r, s) \# ws, N) = & \\
\text{let} & \\
\text{succs} = \text{map } (\lambda b. & \\
\text{let} & \\
r' = \langle \mathcal{D}_b(r) \rangle; \quad s' = \langle \mathcal{D}_b(s) \rangle & \\
\text{in } ((r', s'), (\text{norm } r', \text{norm } s')) \in \Sigma_n; & \\
\text{new} = \text{remdups snd } (\text{filter } (\lambda(-, rs). rs \notin N) \text{succs}) & \\
\text{in } (ws \text{ @ } \text{map fst new}, \text{set } (\text{map snd new}) \cup N) &
\end{aligned}$$

The function $\text{set} :: \alpha \text{ list} \rightarrow \alpha \text{ set}$ maps a list to the set of its elements, $\text{filter} :: (\alpha \rightarrow \mathbf{B}) \rightarrow \alpha \text{ list} \rightarrow \alpha \text{ list}$ removes elements that do not fulfill the given predicate, while $\text{remdups} :: (\alpha \rightarrow \beta) \rightarrow \alpha \text{ list} \rightarrow \alpha \text{ list}$ is used to keep the worklist as small as possible. $\text{remdups } f \text{ xs}$ removes duplicates from xs modulo the function f , e.g., $\text{remdups snd } [(\mathbf{0}, \mathbf{0}), (\mathbf{1}, \mathbf{0})] = [(\mathbf{1}, \mathbf{0})]$ (which element is actually kept is irrelevant; the result $[(\mathbf{0}, \mathbf{0})]$ would also be valid).

Finally, a well-formedness check completes the now executable algorithm $\text{eqv}^{\text{RE}} :: \mathbb{N} \rightarrow \alpha \text{ RE} \rightarrow \alpha \text{ RE} \rightarrow \mathbf{B}$.

$$\begin{aligned}
\text{eqv}_n^{\text{RE}} r s = & \\
& \text{wf}_n(r) \wedge \text{wf}_n(s) \wedge \\
& (\text{case while b } c_n([\langle r \rangle, \langle s \rangle], \{(\text{norm } \langle r \rangle, \text{norm } \langle s \rangle)\}) \text{ of} \\
& \quad \text{Some } ([], _) \Rightarrow \top \\
& \quad | \text{Some } (_ \# _, _) \Rightarrow \perp)
\end{aligned}$$

The termination of eqv^{RE} for any input is guaranteed by two facts: (1) all recursively defined functions in Isabelle/HOL terminate by their definitional principle (either primitive or wellfounded recursion) and (2) the termination of while follows from Theorem 2 and the fact that the set N of already seen pairs in the state is a subset of $(\lambda(r, s). (\text{norm } r, \text{norm } s)) \cdot \{(\langle \mathcal{D}_w^*(r) \rangle, \langle \mathcal{D}_w^*(s) \rangle) \mid w \in \Sigma_n^*\}$.

Theorem 6 (Termination)

Let r and s be n -well-formed regular expressions. Then

$$\text{while b } c_n([\langle r \rangle, \langle s \rangle], \{(\text{norm } \langle r \rangle, \text{norm } \langle s \rangle)\}) \neq \text{None}.$$

Function eqv^{RE} deserves the name decision procedure since it constitutes a refinement of the algorithm abstractly stated in Theorem 5, and is therefore sound and complete. The refinement follows from proving the following predicate being an invariant for the states (ws, N) of the while-loop given two initial n -well-formed regular expressions r and s :

$$\begin{aligned}
\text{inv } (ws, N) = & \\
& (\forall(r', s') \in \text{set } ws. (\text{norm } r', \text{norm } s') \in N) \wedge \\
& (\forall(r', s') \in N. \exists w \in \Sigma_n^*. \mathcal{D}_w^*(r) = r' \wedge \mathcal{D}_w^*(s) = s') \wedge \\
& (\forall(r', s') \in N \setminus ((\lambda(r, s). (\text{norm } r, \text{norm } s)) \cdot (\text{set } ws)). \varepsilon(r') \leftrightarrow \varepsilon(s') \wedge \\
& \quad (\forall a \in \Sigma_n. (\text{norm } (\mathcal{D}_a(r')), \text{norm } (\mathcal{D}_a(s'))) \in N)
\end{aligned}$$

For an execution of eqv^{RE} , either ws is eventually emptied—in which case the last conjunct of inv corresponds to N being a bisimulation modulo norm —or the test

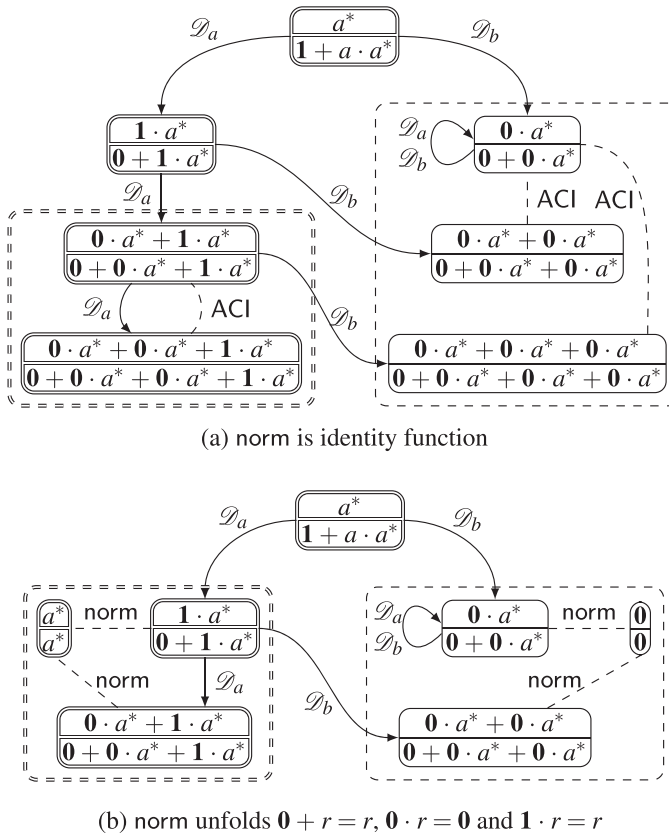


Fig. 1. Checking the equivalence of a^* and $\mathbf{1} + a \cdot a^*$ for $\Sigma_n = \{a, b\}$.

b fails for pair in ws yielding a counterexample to language equivalence using the first two conjuncts of inv .

Theorem 7 (Soundness)

Let r and s be regular expressions such that $eqv_n^{RE} r s$. Then $\mathcal{L}_n(r) = \mathcal{L}_n(s)$.

Theorem 8 (Completeness)

Let r and s be n -well-formed regular expressions such that $\mathcal{L}_n(r) = \mathcal{L}_n(s)$. Then $eqv_n^{RE} r s$.

Let us observe the decision procedure at work by looking at the regular expressions a^* and $\mathbf{1} + a \cdot a^*$ for some $a \in \Sigma_n = \{a, b\}$ for some n . For presentation purposes, the correspondence of derivatives to automata is useful. Figure 1 shows two automata, the states of which are equivalence classes of pairs of regular expressions indicated by a dashed fringe (which is omitted for singleton classes). The equivalence classes of automaton (a) are modulo plain ACI normalization, while those of automaton (b) are modulo a stronger normalization function, making the automaton smaller. Transitions correspond to pairwise derivatives and doubled margins denote states for which the associated pairs of regular expressions are pairwise nullable. Both automata are the result of our decision procedure performing a breadth-first

exploration starting with the initially given pair and ignoring states that are in the equivalence class of already visited states. The absence of pairs (r, s) for which r is nullable and s is not nullable (or vice versa) proves the equivalence of all pairs in the automaton, including the pair $(a^*, \mathbf{1} + a \cdot a^*)$.

Let us mention two obvious performance deficits of our algorithm. First, the ACI normalization is sorting the summands in expressions basically using bubble-sort. Using a set data structure instead of binary sums would improve this to merge-sort and is certainly desirable. Second, the algorithm constructs a bisimulation (not even a bisimulation up to equality). This effectively means that even when applied on two identical expressions, the algorithm would still enumerate all derivatives. There is a whole hierarchy of possible improvements: bisimulation up to equality, equivalence, congruence, and congruence and context, which have been successfully employed in unverified derivative-based decision procedures (Bonchi & Pous, 2013; Pous, 2015). However, when verifying an algorithm one has to settle for a solution somewhere in between of efficiency and simplicity.

4.4 Atoms with more structure

Owens *et al.* (2009) advocate a more compact regular expression structure where the language of an atom denotes a set of one letter words. The gained compactness is beneficial especially for expressions over a large alphabet. In our setting, this would mean using the type $(\alpha \text{ set})$ RE instead of α RE (without changing the underlying alphabet type α). We will see later that our alphabet is indeed large—exponential in the number of free variables.

We generalize this idea without committing to a fixed type for the atoms yet. Instead of α RE, the regular expressions over the alphabet type α on which the algorithm operates will be of type β RE, where the relationship between α and the new atoms β is given by a function $\text{mem}^\wedge :: \beta \rightarrow \alpha \rightarrow \mathbf{B}$. The new semantics $\mathcal{L} :: \mathbf{N} \rightarrow \beta \text{ RE} \rightarrow (\alpha \text{ list}) \text{ set}$ of such regular expressions is defined just as the old \mathcal{L} except for the atom case. A similar adjustment is required for the new derivative $\mathcal{D} :: \alpha \rightarrow \beta \text{ RE} \rightarrow \beta \text{ RE}$.

$$\mathcal{L}_n(b) = \{a \mid \text{mem}^\wedge b a\} \quad \mathcal{D}_a(b) = \text{if mem}^\wedge b a \text{ then } \mathbf{1} \text{ else } \mathbf{0}$$

Furthermore, the function $\text{wf}^\wedge :: \text{nat} \rightarrow \beta \rightarrow \mathbf{B}$ is used to detect whether a β -atom is wellformed. The wellformedness check for regular expressions $\text{wf} :: \text{nat} \rightarrow \beta \text{ RE} \rightarrow \mathbf{B}$ will use wf^\wedge in the atom case: $\text{wf}_n(b) = \text{wf}^\wedge n b$. The functions mem^\wedge and wf^\wedge are two further parameters of our procedure. We obtain the original procedure by instantiating β with α and defining $\text{mem}^\wedge (b :: \alpha) a \leftrightarrow (a = b)$ and $\text{wf}^\wedge n b \leftrightarrow (b \in \Sigma_n)$. For the data structure from Owens *et al.* (2009), one would instantiate β with $\alpha \text{ set}$ and define $\text{mem}^\wedge (B :: \alpha \text{ set}) a \leftrightarrow (a \in B)$ and $\text{wf}^\wedge n B \leftrightarrow (\forall a \in B. a \in \Sigma_n)$.

The benefit of the abstract formulation is the fact that β can be instantiated with a set representation tailored to the particularities of the used regular expressions. In our case, the regular expressions are translated MSO formulas and a few very particular sets of letters arise from the translation. Therefore, in Section 5 we will

define a data type α atom matching exactly those particularities and instantiate β , mem^\wedge , and wf^\wedge accordingly.

4.5 Alternatives to Brzozowski derivatives

The example from Figure 1 shows that the choice of the normalization is crucial for the size of the bisimulation relation. In prior work (Nipkow & Traytel, 2014), we show that partial derivatives² of ordinary regular expressions can be represented by a composition $\langle\langle - \rangle\rangle \circ \mathcal{D}_a$ where $\langle\langle - \rangle\rangle$ is a particular normalization function defined using smart constructors and observe that $\langle\langle - \rangle\rangle$ tends to maintain a better balance between the size of the resulting bisimulation and the ease to compute the normal form than other ad hoc choices. To use partial derivatives here, we extend this particular function $\langle\langle - \rangle\rangle$ to Π -extended regular expressions as follows. The equations for the smart constructors \boxplus , \boxminus , \boxdot , \boxsquare , and \boxplus are matched sequentially.

$$\begin{array}{ll}
\langle\langle \mathbf{0} \rangle\rangle & = \mathbf{0} & \mathbf{0} \boxplus r & = r \\
\langle\langle \mathbf{1} \rangle\rangle & = \mathbf{1} & r \boxplus \mathbf{0} & = r \\
\langle\langle a \rangle\rangle & = a & (r + s) \boxplus t & = r \boxplus (s \boxplus t) \\
\langle\langle r + s \rangle\rangle & = \langle\langle r \rangle\rangle \boxplus \langle\langle s \rangle\rangle & r \boxplus (s + t) & = \text{if } r = s \text{ then } s + t \\
\langle\langle r \cdot s \rangle\rangle & = \langle\langle r \rangle\rangle \boxdot s & & \text{else if } r \leq s \text{ then } r + (s + t) \\
\langle\langle r^* \rangle\rangle & = r^* & & \text{else } s + (r \boxplus t) \\
\langle\langle r \cap s \rangle\rangle & = \langle\langle r \rangle\rangle \boxsquare \langle\langle s \rangle\rangle & r \boxplus s & = \text{if } r = s \text{ then } r \\
\langle\langle \neg r \rangle\rangle & = \boxminus \langle\langle r \rangle\rangle & & \text{else if } r \leq s \text{ then } r + s \\
\langle\langle \Pi r \rangle\rangle & = \boxplus \langle\langle r \rangle\rangle & & \text{else } s + r \\
\mathbf{0} \boxminus r & = \mathbf{0} & \mathbf{0} \boxdot r & = \mathbf{0} \\
\mathbf{1} \boxminus r & = r & r \boxdot \mathbf{0} & = \mathbf{0} \\
(r + s) \boxminus t & = (r \boxminus s) \boxplus (s \boxminus t) & (\neg \mathbf{0}) \boxdot r & = r \\
r \boxminus s & = s \cdot t & r \boxdot (\neg \mathbf{0}) & = r \\
\boxminus (r + s) & = (\boxminus r) \boxdot (\boxminus s) & (r + s) \boxdot t & = (r \boxdot t) \boxplus (s \boxdot t) \\
\boxminus (r \cap s) & = (\boxminus r) \boxplus (\boxminus s) & r \boxdot (s + t) & = (r \boxdot s) \boxplus (r \boxdot t) \\
\boxminus (\neg r) & = r & (r \cap s) \boxdot t & = r \boxdot (s \boxdot t) \\
\boxminus r & = \neg r & r \boxdot (s \cap t) & = \text{if } r = s \text{ then } s \cap t \\
& & & \text{else if } r \leq s \text{ then } r \cap (s \cap t) \\
& & & \text{else } s \cap (r \boxdot t) \\
\boxplus \mathbf{0} & = \mathbf{0} & r \boxdot s & = \text{if } r = s \text{ then } r \\
\boxplus \mathbf{1} & = \mathbf{1} & & \text{else if } r \leq s \text{ then } r \cap s \\
\boxplus (r + s) & = (\boxplus r) \boxplus (\boxplus s) & & \text{else } s \cap r \\
\boxplus r & = \Pi r & &
\end{array}$$

It is worth noticing that $\langle\langle - \rangle\rangle$ does not descend recursively into right-hand side of concatenation and into iteration. Also, \boxsquare distributes over \boxplus , which establishes something like a disjunctive normal form with respect to intersection (conjunction) and union (disjunction). Our motivation for this design goes back to

² Partial derivatives (Antimirov, 1996) refine Brzozowski derivatives by splitting the derivation result at some $+$ -constructors into a finite set of regular expressions. Partial derivatives correspond to nondeterministic automata in the same way derivatives correspond to deterministic ones.

Caron *et al.* (2011), who show how to extend partial derivatives to negation and intersection using sets of sets of regular expressions. The outer level of sets there represents unions, the inner intersections. We conjecture that the usage of our $\langle\langle - \rangle\rangle$ as the normalization function produces isomorphic bisimulations to those obtained by working with the extended partial derivatives by Caron *et al.* (2011) directly, but do not attempt to prove it. This conjecture is irrelevant for our purpose, since there is anyway only empirical evidence that partial derivatives perform better than other normalizations for Π -extended regular expression, yet it is an interesting problem to work on in the future. To employ $\langle\langle - \rangle\rangle$ in the algorithm, it is sufficient to prove that it preserves wellformedness and languages—an easy exercise in induction.

Lemma 9

Let r be an n -well-formed regular expression. Then $\text{wf}_n\langle\langle r \rangle\rangle$ and $\mathcal{L}_n\langle\langle r \rangle\rangle = \mathcal{L}_n(r)$.

We remark that the normalization $\langle\langle - \rangle\rangle$ does not enjoy nice algebraic properties. The source of the problem is that our smart constructor \square is not idempotent. To see this, assuming $a \leq b \leq a \cap b$, we calculate: $(a + b) \square (a + b) = a + b + (a \cap b)$. Consequently, the de Morgan law $\langle\langle \neg(r + s) \rangle\rangle = \langle\langle \neg r \cap \neg s \rangle\rangle$ does not hold. One could argue that this is a bad design of the normalization, which is modeled after the operations on sets of sets of expressions given elsewhere (Caron *et al.*, 2011). (Those operations suffer from the same limitations.) However, the performance when using this normalization in practice seems reasonable and our attempts in changing the normalization function to make \square idempotent (for example, by giving up distributivity of \cap over $+$ or by adding more equality checks in the definition of \square) resulted in a perceivable decrease in performance.

Nevertheless, an interesting question is whether one can find a fast normalization function that decides equivalence under the following inductively defined equivalence relation \approx , which is modeled after what the normalization function $\langle\langle - \rangle\rangle$ attempts (but fails) to equate. Note that, unlike \sim (ACI), the relation \approx is only an equivalence, not a congruence. Not being able to find such a normalization, we leave this question as future work.

$$\begin{array}{llll}
 \mathbf{0} + r \approx r & r + \mathbf{0} \approx r & \mathbf{0} \cdot r \approx \mathbf{0} & \mathbf{1} \cdot r \approx r \\
 (\neg \mathbf{0}) \cap r \approx r & r \cap (\neg \mathbf{0}) \approx r & \mathbf{0} \cap r \approx \mathbf{0} & r \cap \mathbf{0} \approx \mathbf{0} \\
 r + (s + t) \approx (r + s) + t & r + s \approx s + r & r + r \approx r & \\
 r \cap (s \cap t) \approx (r \cap s) \cap t & r \cap s \approx s \cap r & r \cap r \approx r & \\
 r \cap (s + t) \approx (r \cap s) + (r \cap t) & (r + s) \cap t \approx (r \cap t) + (s \cap t) & & \\
 (r + s) \cdot t \approx (r \cdot t) + (s \cdot t) & \neg(\neg r) \approx r & & \\
 \neg(r + s) \approx (\neg r) \cap (\neg s) & \neg(r \cap s) \approx (\neg r) + (\neg s) & & \\
 \Pi(r + s) \approx \Pi r + \Pi s & \Pi \mathbf{0} \approx \mathbf{0} & \Pi \mathbf{1} \approx \mathbf{1} & \\
 r \approx r & \frac{r \approx s}{s \approx r} & \frac{r \approx s \quad s \approx t}{r \approx t} & \frac{r_1 \approx s_1 \quad r_2 \approx s_2}{r_1 + r_2 \approx s_1 + s_2} \\
 \frac{r_1 \approx s_1 \quad r_2 \approx s_2}{r_1 \cap r_2 \approx s_1 \cap s_2} & \frac{r_1 \approx s_1}{r_1 \cdot t \approx s_1 \cdot t} & \frac{r \approx s}{\neg r \approx \neg s} & \frac{r \approx s}{\Pi r \approx \Pi s}
 \end{array}$$

Another promising alternative to Brzozowski derivatives is the data type $\alpha \text{RE}^{\text{op}}$ of dual regular expressions (Okhotin, 2005). The data type is obtained by modifying

α RE as following: drop the negation and intersection constructors and add to every remaining n -ary constructor \circ a Boolean flag b with the following semantics $\mathcal{L}^{\text{op}} :: \mathbb{N} \rightarrow \alpha \text{ RE}^{\text{op}} \rightarrow \mathbb{B}$.

$$\mathcal{L}_n^{\text{op}}(\circ b r_1 \cdots r_n) = \mathcal{L}_n(\text{if } b \text{ then } \circ r_1 \cdots r_n \text{ else } \neg(\circ(\neg r_1) \cdots (\neg r_n)))$$

In the formalization (Traytel & Nipkow, 2014) we define wellformedness, derivatives and some ad hoc normalization on $\alpha \text{ RE}^{\text{op}}$ and generalize the bisimulation construction to work on $\alpha \text{ RE}^{\text{op}}$ as well. The evaluation will show that the decision procedure we obtain by using dual regular expression performs better for the WS1S semantics of MSO, but worse for the M2L semantics—a phenomenon for which we do not have an explanation yet.

5 MSO on finite words

Logics on finite words consider formulas in the context of a formal word, with variables representing positions in the word. In the first-order logic on words, a variable always denotes a single position while in MSO logic on finite words, variables come in two flavors: first-order variables for single positions and second-order variables for finite sets of positions.

In the next subsections, we first define the syntax of formulas and give them a semantics that is related to formal languages: M2L(Str). The second semantics, WS1S, is introduced as a relaxation of M2L (we drop the “(Str)” from now on). Both semantics are equally expressive and deciding both is of nonelementary complexity. The benefits and drawbacks of the two semantics are discussed elsewhere (Klarlund, 1999; Ayari & Basin, 2000).

5.1 Syntax and M2L semantics

MSO formulas are syntactically represented by the recursive data type $\alpha \Phi$ using de Bruijn indices for variable bindings. Terms of $\alpha \Phi$ are generated by the grammar

$$\varphi = \mathbf{Q} a m \mid m_1 < m_2 \mid m \in M \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \neg \varphi \mid \exists \varphi \mid \exists \varphi$$

where $\varphi, \psi :: \alpha \Phi$, $m, m_1, m_2, M :: \mathbb{N}$ and $a :: \alpha$. Lower-case variables m, m_1, m_2 denote first-order variables, M denotes a second-order variable. The atomic formula $\mathbf{Q} a m$ requires the letter of the word at the position represented by variable m to be a ; the constructors $<$ and \in compare positions; Boolean operators are interpreted as usual.

The bold existential quantifier \exists binds second-order variables, \exists binds first-order variables. Occurrences of bound variables represented as de Bruijn indices refer to their binders by counting the number of nested existential quantifier between the binder and the occurrence. For example, the formula $\exists(\mathbf{Q} a 0 \wedge (\exists 1 \in 0))$ translates to $\exists x.(\mathbf{Q} a x \wedge (\exists X. x \in X))$ when using names. The first 0 in the nameless formula refers to the outermost first-order quantifier. Inside of the inner second-order quantifier, index 1 refers to the outermost quantifier and index 0 to the inner quantifier. The nameless representation simplifies reasoning by implicitly capturing

α -equivalence of formulas. On the downside, de Bruijn indices are less readable and must be manipulated with care.

Formulas may have free variables. The functions $\mathcal{V}_1 :: \alpha \Phi \rightarrow \mathbb{N}$ set and $\mathcal{V}_2 :: \alpha \Phi \rightarrow \mathbb{N}$ set collect the free first-order and second-order variables:

$$\begin{array}{ll}
 \mathcal{V}_1(Q a m) = \{m\} & \mathcal{V}_2(Q a m) = \{\} \\
 \mathcal{V}_1(m_1 < m_2) = \{m_1, m_2\} & \mathcal{V}_2(m_1 < m_2) = \{\} \\
 \mathcal{V}_1(m \in M) = \{m\} & \mathcal{V}_2(m \in M) = \{M\} \\
 \mathcal{V}_1(\varphi \wedge \psi) = \mathcal{V}_1(\varphi) \cup \mathcal{V}_1(\psi) & \mathcal{V}_2(\varphi \wedge \psi) = \mathcal{V}_2(\varphi) \cup \mathcal{V}_2(\psi) \\
 \mathcal{V}_1(\varphi \vee \psi) = \mathcal{V}_1(\varphi) \cup \mathcal{V}_1(\psi) & \mathcal{V}_2(\varphi \vee \psi) = \mathcal{V}_2(\varphi) \cup \mathcal{V}_2(\psi) \\
 \mathcal{V}_1(\neg \varphi) = \mathcal{V}_1(\varphi) & \mathcal{V}_2(\neg \varphi) = \mathcal{V}_2(\varphi) \\
 \mathcal{V}_1(\exists \varphi) = \lfloor \mathcal{V}_1(\varphi) \setminus \{0\} \rfloor & \mathcal{V}_2(\exists \varphi) = \lfloor \mathcal{V}_2(\varphi) \rfloor \\
 \mathcal{V}_1(\exists \varphi) = \lfloor \mathcal{V}_1(\varphi) \rfloor & \mathcal{V}_2(\exists \varphi) = \lfloor \mathcal{V}_2(\varphi) \setminus \{0\} \rfloor
 \end{array}$$

The notation $\lfloor X \rfloor$ is shorthand for $(\lambda x. x - 1) \cdot X$, which reverts the increasing effect of an existential quantifier on previously bound or free variables. To obtain only free variables, bound variables are removed when their quantifier is processed, at which point the bound variable has index 0.

Just as for Π -extended regular expressions, not all formulas in $\alpha \Phi$ are meaningful. Consider $0 \in 0$, where 0 is both a first-order and a second-order variable. To exclude such formulas, we define the predicate $\text{wf}^\Phi :: \mathbb{N} \rightarrow \alpha \Phi \rightarrow \mathbb{B}$ as $\text{wf}_n^\Phi(\varphi) = (\mathcal{V}_1(\varphi) \cap \mathcal{V}_2(\varphi) = \{\}) \wedge \text{pre_wf}_n^\Phi(\varphi)$ and call a formula φ *n-wellformed* if $\text{wf}_n^\Phi(\varphi)$ holds. The recursively defined predicate $\text{pre_wf}^\Phi :: \mathbb{N} \rightarrow \alpha \Phi \rightarrow \mathbb{B}$ is used for further assumptions on the structure of *n-well-formed* formulas, which will simplify our proofs:

$$\begin{array}{ll}
 \text{pre_wf}_n^\Phi(Q a m) = a \in \Sigma \wedge m < n & \\
 \text{pre_wf}_n^\Phi(m_1 < m_2) = m_1 < n \wedge m_2 < n & \\
 \text{pre_wf}_n^\Phi(m \in M) = m < n \wedge M < n & \\
 \text{pre_wf}_n^\Phi(\varphi \wedge \psi) = \text{pre_wf}_n^\Phi(\varphi) \wedge \text{pre_wf}_n^\Phi(\psi) & \\
 \text{pre_wf}_n^\Phi(\varphi \vee \psi) = \text{pre_wf}_n^\Phi(\varphi) \wedge \text{pre_wf}_n^\Phi(\psi) & \\
 \text{pre_wf}_n^\Phi(\neg \varphi) = \text{pre_wf}_n^\Phi(\varphi) & \\
 \text{pre_wf}_n^\Phi(\exists \varphi) = \text{pre_wf}_{n+1}^\Phi(\varphi) \wedge 0 \in \mathcal{V}_1(\varphi) \wedge 0 \notin \mathcal{V}_2(\varphi) & \\
 \text{pre_wf}_n^\Phi(\exists \varphi) = \text{pre_wf}_{n+1}^\Phi(\varphi) \wedge 0 \notin \mathcal{V}_1(\varphi) \wedge 0 \in \mathcal{V}_2(\varphi) &
 \end{array}$$

$\text{pre_wf}_n^\Phi(\varphi)$ ensures that the index of every free variable in φ is below n and the values of type α come from a fixed alphabet Σ . Note that Σ is really just a fixed set of letters of type α , independent of any n and is a parameter of our setup. Moreover, pre_wf^Φ checks that bound variables are correctly used as first-order or second-order with respect to their binders and excludes formulas with unused binders; unused binders are obviously superfluous.

An *interpretation* of an MSO formula is a pair of a word $w :: \alpha$ list from Σ^* and an assignment $\mathcal{I} :: (\mathbb{N} + \mathbb{N}$ set) list for free variables. The latter essentially consists of two functions with finite domain: one from first-order variables to positions and the other from second-order variables to sets of positions. We represent those two functions by a list, once again benefiting from de Bruijn indices—the value lookup for a variable with de Bruijn index i corresponds to inspecting the assignment \mathcal{I} at

position i , i.e., $\mathcal{I}[i]$. The range of \mathcal{I} is a sum type, denoting the disjoint union of its two argument types. The sum type has two constructors $\text{Inl} :: \alpha \rightarrow \alpha + \beta$ and $\text{Inr} :: \beta \rightarrow \alpha + \beta$, such that for a first-order variable m there is a position p with $\mathcal{I}[m] = \text{Inl } p$ and for a second-order variable M there is a finite set of positions P with $\mathcal{I}[M] = \text{Inr } P$.

An interpretation that *satisfies* a formula is called a model. Satisfiability for M2L, denoted by infix $\models :: \alpha \text{ list} \times (\mathbb{N} + \mathbb{N} \text{ set}) \text{ list} \rightarrow \alpha \Phi \rightarrow \mathbb{B}$, is defined recursively on $\alpha \Phi$. To simplify the notation, the constructors Inl and Inr are stripped implicitly in the definition.

$$\begin{aligned}
(w, \mathcal{I}) \models \mathbf{Q} a m &\leftrightarrow w[\mathcal{I}[m]] = a \\
(w, \mathcal{I}) \models m_1 < m_2 &\leftrightarrow \mathcal{I}[m_1] < \mathcal{I}[m_2] \\
(w, \mathcal{I}) \models m \in M &\leftrightarrow \mathcal{I}[m] \in \mathcal{I}[M] \\
(w, \mathcal{I}) \models \varphi \wedge \psi &\leftrightarrow (w, \mathcal{I}) \models \varphi \wedge (w, \mathcal{I}) \models \psi \\
(w, \mathcal{I}) \models \varphi \vee \psi &\leftrightarrow (w, \mathcal{I}) \models \varphi \vee (w, \mathcal{I}) \models \psi \\
(w, \mathcal{I}) \models \neg \varphi &\leftrightarrow (w, \mathcal{I}) \not\models \varphi \\
(w, \mathcal{I}) \models \exists \varphi &\leftrightarrow \exists p \in \{0, \dots, |w| - 1\}. (w, \text{Inl } p \# \mathcal{I}) \models \varphi \\
(w, \mathcal{I}) \models \exists \varphi &\leftrightarrow \exists P \subseteq \{0, \dots, |w| - 1\}. (w, \text{Inr } P \# \mathcal{I}) \models \varphi
\end{aligned}$$

For the definition to make sense, \mathcal{I} must correctly map first-order variables to positions (i.e., $\mathcal{I}[m] = \text{Inl } p$) and second-order variables to sets of positions (i.e., $\mathcal{I}[M] = \text{Inr } P$). Furthermore, all positions in \mathcal{I} should be below the length of the word, and for technical reasons the word should not be empty. We formalize these assumptions by the predicate $\text{wf}^{\text{M2L}} :: \alpha \Phi \rightarrow \alpha \text{ list} \times (\mathbb{N} + \mathbb{N} \text{ set}) \text{ list} \rightarrow \mathbb{B}$ and call an interpretation *M2L-wellformed* for φ if $\text{wf}_\varphi^{\text{M2L}}(w, \mathcal{I})$ holds:

$$\begin{aligned}
\text{wf}_\varphi^{\text{M2L}}(w, \mathcal{I}) = & w \neq [] \wedge w \in \Sigma^* \wedge \\
& \forall \text{Inl } p \in \text{set } \mathcal{I}. p < |w| \wedge \\
& \forall \text{Inr } P \in \text{set } \mathcal{I}. (\forall p \in P. p < |w|) \wedge \\
& \forall m \in \mathcal{V}_1(\varphi). (\exists p. \mathcal{I}[m] = \text{Inl } p) \wedge \\
& \forall M \in \mathcal{V}_2(\varphi). (\exists P. \mathcal{I}[M] = \text{Inr } P)
\end{aligned}$$

5.2 WS1S semantics

In an M2L-well-formed model, positions are restricted by the length of the word. This is the key difference compared to WS1S. In WS1S, no *a priori* restrictions on the variable ranges are made, although all second-order variables still represent finite sets. The subtle difference is illustrated by the formula $\exists (\forall 0 \in 1)$ (with names: $\exists X. \forall x. x \in X$), where $\forall \varphi$ is just an abbreviation for $\neg \exists \neg \varphi$. In the M2L semantics, $\exists (\forall 0 \in 1)$ is satisfied by all well-formed interpretations—the witness set for the outer existential quantifier is for a well-formed interpretation (w, \mathcal{I}) just the set $\{0, \dots, |w| - 1\}$. In contrast, in WS1S, there is no finite set which contains all arbitrarily large positions, thus $\exists (\forall 0 \in 1)$ is unsatisfiable.

Formally, satisfiability for WS1S, denoted by infix $\models :: \alpha \text{ list} \times (\mathbb{N} + \mathbb{N} \text{ set}) \text{ list} \rightarrow \alpha \Phi \rightarrow \mathbb{B}$, is defined just as for M2L (replacing \models by \models) except for the following

equations.

$$\begin{aligned}
 (w, \mathcal{I}) \models \mathbf{Q} a m &\leftrightarrow (\text{if } \mathcal{I}[m] < |w| \text{ then } w[\mathcal{I}[m]] \text{ else } z) = a \\
 (w, \mathcal{I}) \models \exists \varphi &\leftrightarrow \exists p. (w, \text{Inl } p \# \mathcal{I}) \models \varphi \\
 (w, \mathcal{I}) \models \exists \varphi &\leftrightarrow \exists P. (w, \text{Inr } P \# \mathcal{I}) \models \varphi \wedge \text{finite } P
 \end{aligned}$$

Here, z is a distinguished letter from Σ . WS1S as defined in the literature does not handle the $\mathbf{Q} a m$ case at all, usually interpreting formulas only with respect to the assignment \mathcal{I} . In order to be able to use the same syntax and the same type of interpretations for both semantics, we have made the above choice. This also allows us to translate $\mathbf{Q} a m$ into the same regular expression irrespective of the intended semantics.

Besides the mentioned relaxation of *WS1S-wellformedness* regarding variable ranges, the empty word also does not impose technical complications as in M2L. Therefore, the predicate $\text{wf}^{\text{WS1S}} :: \alpha \Phi \rightarrow \alpha \text{list} \times (\mathbb{N} + \mathbb{N} \text{ set}) \text{list} \rightarrow \mathbb{B}$ is defined as follows.

$$\begin{aligned}
 \text{wf}_{\varphi}^{\text{WS1S}}(w, \mathcal{I}) &= w \in \Sigma^* \wedge \\
 &\quad \forall \text{Inr } P \in \text{set } \mathcal{I}. \text{finite } P \wedge \\
 &\quad \forall m \in \mathcal{V}_1(\varphi). (\exists p. \mathcal{I}[m] = \text{Inl } p) \wedge \\
 &\quad \forall M \in \mathcal{V}_2(\varphi). (\exists P. \mathcal{I}[M] = \text{Inr } P)
 \end{aligned}$$

5.3 Encoding interpretations as words

Formulas are equivalent if they have the same set of well-formed models. To relate equivalent formulas with language equivalent regular expressions, the set of well-formed models must be represented as a formal language by encoding interpretations as words. As before, we cover the encoding of the M2L semantics first.

To simplify the formalization, we choose a very simple encoding using Boolean vectors. For an interpretation (w, \mathcal{I}) , we associate with every position p in the word w a Boolean vector bs of length $|\mathcal{I}|$, such that $bs[m] = \top$ iff the m th variable in \mathcal{I} is first-order and its value is p or it is second-order and its value contains p . For example, for $\Sigma = \{a, b\}$ the interpretation $(w, \mathcal{I}) = (aba, \text{Inl } 0 \# \text{Inr } \{1, 2\} \# \text{Inl } 2 \# \square)$ can be written in two dimensions as follows:

	a	b	a
Inl 0	⊤	⊥	⊥
Inr {1, 2}	⊥	⊤	⊤
Inl 2	⊥	⊥	⊤

In the first row, the value \top is placed only in the first column because the first variable of \mathcal{I} is the first-order position 0. In general, the columns correspond to the Boolean vectors associated with positions in the word, while every row corresponds to one variable. For first-order variables, there must be exactly one \top per row. The first row encodes the value of the most recently bound \mathcal{I} variable. Now, we consider every column as a letter of a new alphabet, which is the underlying alphabet $\Sigma_n = \Sigma \times \mathbb{B}^n$ of regular expressions of Section 4. This transformation of interpretations into

words over Σ_n is performed by the function $\text{enc}^{\text{M2L}} :: \alpha \text{ list} \times (\mathbb{N} + \mathbb{N} \text{ set}) \text{ list} \rightarrow (\alpha \times \mathbb{B} \text{ list}) \text{ list}$; we omit its obvious definition.

Furthermore, the second parameter $\pi :: \Sigma_{n+1} \rightarrow \Sigma_n$ of our decision procedure for regular expressions can now be instantiated as the function that maps $(a, b \# bs)$ to (a, bs) . Thus, the projection Π operates on words by removing the first row from words in the language of the body expression, reflecting the semantics of an existential quantifier.

Finally, the *M2L-language* $\mathcal{L}^{\text{M2L}} :: \mathbb{N} \rightarrow \alpha \Phi \rightarrow (\alpha \times \mathbb{B} \text{ list}) \text{ set}$ of an MSO formula is the set of encodings of its well-formed models, i.e., $\mathcal{L}_n^{\text{M2L}}(\varphi) = \{\text{enc}^{\text{M2L}}(w, \mathcal{I}) \mid \text{wf}_\varphi^{\text{M2L}}(w, \mathcal{I}) \wedge |\mathcal{I}| = n \wedge (w, \mathcal{I}) \models \varphi\}$.

Concerning WS1S, the encoding is slightly more complicated due to the following observation: Interpretations (w, \mathcal{I}) and (wz^n, \mathcal{I}) for all $n :: \mathbb{N}$ behave the same when considering satisfiability and wellformedness with respect to a formula (z^n denotes n -fold repetition of the letter z as a word). That suggests that the example interpretation $(w, \mathcal{I}) = (aba, \text{Inl } 0 \# \text{Inr } \{1, 2\} \# \text{Inl } 2 \# \square)$ from above can be encoded as

	a	b	a	z^m
Inl 0	T	⊥	⊥	⊥ ^m
Inr {1, 2}	⊥	T	T	⊥ ^m
Inl 2	⊥	⊥	T	⊥ ^m

for every $m :: \mathbb{N}$. Hence, a single WS1S interpretation is translated into a countably infinite set of words by a function $\text{enc}^{\text{WS1S}} :: \alpha \text{ list} \times (\mathbb{N} + \mathbb{N} \text{ set}) \text{ list} \rightarrow (\alpha \times \mathbb{B} \text{ list}) \text{ list set}$; we again omit its formal definition. Accordingly, the *WS1S-language* $\mathcal{L}^{\text{WS1S}} :: \mathbb{N} \rightarrow \alpha \Phi \rightarrow (\alpha \times \mathbb{B} \text{ list}) \text{ set}$ of an MSO formula is defined as the union of all encodings of its well-formed models: $\mathcal{L}_n^{\text{WS1S}}(\varphi) = \bigcup \{\text{enc}^{\text{WS1S}}(w, \mathcal{I}) \mid \text{wf}_\varphi^{\text{WS1S}}(w, \mathcal{I}) \wedge |\mathcal{I}| = n \wedge (w, \mathcal{I}) \models \varphi\}$.

5.4 From M2L formulas to regular expressions

We have fixed the underlying alphabet type $\alpha \times \mathbb{B} \text{ list}$ of the language of a formula. In principle, we could start translating formulas of type $\alpha \Phi$ into regular expressions of type $(\alpha \times \mathbb{B} \text{ list}) \text{ RE}$. However, the abstraction for atoms introduced in Section 4.4 caters for a more efficient encoding of formulas. We define the data type $\alpha \text{ atom}$ as

$$\text{atom} = \text{A } a_bs \mid \text{AQ } m \ a \mid \text{ANth } m \ b \mid \text{Anth}_2 \ m \ M$$

where $\text{atom} :: \alpha \text{ atom}$, $a_bs :: \alpha \times \mathbb{B} \text{ list}$, $m, m_1, m_2, M :: \mathbb{N}$, $a :: \alpha$, and $b :: \mathbb{B}$. Each constructor of $\alpha \text{ atom}$ represents a set of elements of type $\alpha \times \mathbb{B} \text{ list}$. The constructor **A** represents the singleton set containing the constructor's argument, **AQ** m a encodes all pairs whose first element is a and whose second element (a Boolean vector) has \top at index m . Both, this informal description as well as the constructor name should indicate that **AQ** m a is closely related to the formula **Q** a m . The remaining two constructors have a similar purpose, being related to the other base cases of the formula type. Let us make this precise by instantiating the two

parameters $\text{mem}^\wedge :: \alpha \text{ atom} \rightarrow \alpha \times \mathbb{B} \text{ list} \rightarrow \mathbb{B}$ and $\text{wf}^\wedge :: \mathbb{N} \rightarrow \alpha \text{ atom} \rightarrow \mathbb{B}$.

$$\begin{aligned} \text{mem}^\wedge (A \ a _bs) \ a_bs' &\leftrightarrow a_bs = a_bs' & \text{wf}^\wedge \ n \ (A \ a _bs) &\leftrightarrow a_bs \in \Sigma_n \\ \text{mem}^\wedge (AQ \ m \ a) \ (a', \ bs) &\leftrightarrow a = a' \wedge bs[m] & \text{wf}^\wedge \ n \ (AQ \ m \ a) &\leftrightarrow a \in \Sigma \wedge m < n \\ \text{mem}^\wedge (ANth \ m \ b) \ (_, \ bs) &\leftrightarrow bs[m] = b & \text{wf}^\wedge \ n \ (ANth \ m \ b) &\leftrightarrow m < n \\ \text{mem}^\wedge (ANth_2 \ m \ M) \ (_, \ bs) &\leftrightarrow bs[m] \wedge bs[M] & \text{wf}^\wedge \ n \ (ANth_2 \ m \ M) &\leftrightarrow m < n \wedge M < n \end{aligned}$$

Now, we are set to tackle the translations of formulas into regular expressions. MSO formulas interpreted in M2L are translated by means of the primitive recursive function $\text{mkRE}^{\text{M2L}} :: \mathbb{N} \rightarrow \alpha \Phi \rightarrow (\alpha \text{ atom}) \text{ RE}$.

$$\begin{aligned} \text{mkRE}_n^{\text{M2L}}(Q \ a \ m) &= \neg \mathbf{0} \cdot AQ \ m \ a \cdot \neg \mathbf{0} \\ \text{mkRE}_n^{\text{M2L}}(m_1 < m_2) &= \neg \mathbf{0} \cdot ANth \ m_1 \top \cdot \neg \mathbf{0} \cdot ANth \ m_2 \top \cdot \neg \mathbf{0} \\ \text{mkRE}_n^{\text{M2L}}(m \in M) &= \neg \mathbf{0} \cdot ANth_2 \ m \ M \cdot \neg \mathbf{0} \\ \text{mkRE}_n^{\text{M2L}}(\varphi \wedge \psi) &= \text{mkRE}_n^{\text{M2L}}(\varphi) \cap \text{mkRE}_n^{\text{M2L}}(\psi) \\ \text{mkRE}_n^{\text{M2L}}(\varphi \vee \psi) &= \text{mkRE}_n^{\text{M2L}}(\varphi) + \text{mkRE}_n^{\text{M2L}}(\psi) \\ \text{mkRE}_n^{\text{M2L}}(\neg \varphi) &= \neg \text{mkRE}_n^{\text{M2L}}(\varphi) \\ \text{mkRE}_n^{\text{M2L}}(\exists \varphi) &= \Pi (\text{mkRE}_{n+1}^{\text{M2L}}(\varphi) \cap \text{WF}_{n+1}\{0\}) \\ \text{mkRE}_n^{\text{M2L}}(\exists \varphi) &= \Pi (\text{mkRE}_{n+1}^{\text{M2L}}(\varphi)) \end{aligned}$$

At first, we ignore the function WF that is used in the case of the first-order quantifier. The natural number parameter of mkRE^{M2L} indicates the number for free variables for the processed formula. The parameter is increased when entering recursively the scope of an existential quantifier.

The intuition behind the translation is demonstrated by the case $Q \ a \ m$. We fix a well-formed model (w, \mathcal{I}) of $Q \ a \ m$. This model must satisfy $w[\mathcal{I}[m]] = a$, or equivalently the fact that there exists a Boolean vector bs of length n such that $\text{enc}^{\text{M2L}}(w, \mathcal{I})[\mathcal{I}[m]] = (a, bs)$ and $bs[m] = \top$. Therefore, the letter at position $\mathcal{I}[m]$ of $\text{enc}^{\text{M2L}}(w, \mathcal{I})$ is matched by the “middle” part $AQ \ m \ a$ of $\text{mkRE}_n^{\text{M2L}}(Q \ a \ m)$, while the subexpressions $\neg \mathbf{0}$ (whose language is Σ_n^*) match the first $\mathcal{I}[m]$ and the last $n - \mathcal{I}[m]$ letters of $\text{enc}^{\text{M2L}}(w, \mathcal{I})$.

Conversely, if we fix a word from $\text{mkRE}_n^{\text{M2L}}(Q \ a \ m)$, it will be equal to an encoding of an interpretation that satisfies $Q \ a \ m$ by a similar argument. However, the interpretation might be not wellformed for $Q \ a \ m$. This happens because the regular expression $\text{mkRE}_n^{\text{M2L}}(Q \ a \ m)$ does not capture the distinction between first-order and second-order variables: it accepts encodings of interpretations that have the value \top more than once at different positions representing the same first-order variable. This indicates that the subexpressions $\neg \mathbf{0}$ in the base cases are not precise enough, but also in the case of Boolean operators similar issues arise. So instead of tinkering with the base cases, it is better to separate the generation a regular expression that encodes models from the one that encodes well-formed interpretations.

To rule out not well-formed interpretations is exactly the purpose of the $\text{WF} :: \mathbb{N} \rightarrow \mathbb{N} \text{ set} \rightarrow (\alpha \text{ atom}) \text{ RE}$ function.

$$\text{WF}_n(X) = \bigcap_{m \in X} (ANth \ m \ \perp)^* \cdot ANth \ m \ \top \cdot (ANth \ m \ \perp)^*$$

The regular expression $WF_n(\mathcal{V}_1(\varphi))$ accepts exactly the encodings of well-formed interpretations (both models and non-models) for φ by ensuring that first-order variables are encoded correctly (i.e., forcing the encoding of an interpretation to contain exactly one \top in rows belonging to a first-order variable).

Lemma 10

Let φ be an n -well-formed formula. Then

- $\mathcal{L}_n(WF_n(\mathcal{V}_1(\varphi))) \setminus \{\perp\} = \{\mathbf{enc}^{\text{M2L}}(w, \mathcal{I}) \mid \mathbf{wf}_\varphi^{\text{M2L}}(w, \mathcal{I}) \wedge |\mathcal{I}| = n\}$, and
- $\mathcal{L}_n(WF_n(\mathcal{V}_1(\varphi))) = \{\mathbf{enc}^{\text{WS1S}}(w, \mathcal{I}) \mid \mathbf{wf}_\varphi^{\text{WS1S}}(w, \mathcal{I}) \wedge |\mathcal{I}| = n\}$.

Using WF in every case of the recursive definition of mkRE^{M2L} is sound but very redundant—instead it is enough to perform the intersection once globally for the entire formula and additionally for every variable introduced by the first-order existential quantifier.

MSO formulas interpreted in WS1S are translated into regular expressions by means of the function $\text{mkRE}^{\text{WS1S}} :: \mathbb{N} \rightarrow \alpha \Phi \rightarrow (\alpha \text{ atom}) \text{ RE}$.

The definition of $\text{mkRE}^{\text{WS1S}}$ coincides with the one of mkRE^{M2L} except for the existential quantifier cases:

$$\begin{aligned} \text{mkRE}_n^{\text{WS1S}}(\exists \varphi) &= \mathcal{Q}(z, \perp^n) (\Pi (\text{mkRE}_{n+1}^{\text{WS1S}}(\varphi) \cap WF_{n+1}\{0\})) \\ \text{mkRE}_n^{\text{WS1S}}(\exists \varphi) &= \mathcal{Q}(z, \perp^n) (\Pi (\text{mkRE}_{n+1}^{\text{WS1S}}(\varphi))) \end{aligned}$$

The regular operation $\mathcal{Q} :: \alpha \times \mathbb{B} \text{ list} \rightarrow (\alpha \text{ atom}) \text{ RE} \rightarrow (\alpha \text{ atom}) \text{ RE}$ reestablishes the invariant of having all words terminated with a suffix $(z, \perp^n)^m$ for every $m :: \mathbb{N}$ in the WS1S language encoding of a formula as required by definition of $\mathbf{enc}^{\text{WS1S}}$ (this invariant might be violated by the projection). More precisely, the following language identity holds for an n -well-formed regular expression r :

$$\mathcal{L}_n(\mathcal{Q} a r) = \{xa^m \mid \exists l. xa^l \in \mathcal{L}_n(r)\}$$

The concrete executable definition of \mathcal{Q} is more involved. On a high-level, \mathcal{Q} is computed by repeatedly deriving from the right by a via the function \mathcal{D}_a^- (followed by ACI-normalization) until a repetition is encountered. The definition of \mathcal{D}^- is identical to the familiar \mathcal{D} which derives from the left except for the concatenation and iteration cases (in which it is dual).

$$\begin{aligned} \mathcal{D}_b^-(\mathbf{0}) &= \mathbf{0} & \mathcal{D}_b^-(\mathbf{1}) &= \mathbf{0} \\ \mathcal{D}_b^-(a) &= \text{if } a = b \text{ then } \mathbf{1} \text{ else } \mathbf{0} & \mathcal{D}_b^-(r + s) &= \mathcal{D}_b^-(r) + \mathcal{D}_b^-(s) \\ \mathcal{D}_b^-(r \cdot s) &= & \mathcal{D}_b^-(r^*) &= r^* \cdot \mathcal{D}_b^-(r) \\ &\text{if } \varepsilon(s) \text{ then } r \cdot \mathcal{D}_b^-(s) + \mathcal{D}_b^-(r) & & \\ &\text{else } r \cdot \mathcal{D}_b^-(s) & & \\ \mathcal{D}_b^-(r \cap s) &= \mathcal{D}_b^-(r) \cap \mathcal{D}_b^-(s) & \mathcal{D}_b^-(\neg r) &= \neg \mathcal{D}_b^-(r) \\ \mathcal{D}_b^-(\Pi r) &= \Pi \left(\bigoplus_{c \in \pi^- b} \mathcal{D}_c^-(r) \right) & & \end{aligned}$$

Repeated derivation is implemented using the while combinator. The state over which the combinator iterates is of type $\mathbb{B} \times \alpha \text{ RE list}$. The Boolean component simply indicates whether the loop should be executed once more, while the list

contains all the derivatives computed so far in reversed order (i.e., the last element of the list is the initial regular expression). The loop is exited on the first *déjà vu*. The termination of this procedure is established by the dual of Theorem 2 for ACI-equivalent “right derivatives”. After exiting the while loop, the operation \mathcal{Q} unions all the expressions computed so far, yielding an operation whose language is $\{x \mid \exists l. xa^l \in \mathcal{L}_n(r)\}$. To obtain the desired semantics, the iteration of a (lifted to the α atom type by A) is concatenated to the union.

$$\begin{aligned} \mathbf{b}^- (\text{continue}, _) &= \text{continue} \\ \mathbf{c}_a^- (_, rs) &= \\ &\text{let } s = \langle \mathcal{D}_a^+ (\text{head } rs) \rangle \\ &\text{in if } s \in \text{set } rs \text{ then } (\perp, rs) \text{ else } (\top, s \# rs) \\ \mathcal{Q} a r &= \\ &\text{let } R = \text{case while } \mathbf{b}^- \mathbf{c}_a^- (\top, \langle r \rangle) \text{ of Some } (_, rs) \Rightarrow \text{set } rs \\ &\text{in } \left(\bigoplus_{r \in R} r \right) \cdot (A a)^* \end{aligned}$$

Finally, we can establish the language correspondence between formulas and generated regular expressions.

Theorem 11

Let φ be an n -well-formed formula. Then

- $\mathcal{L}_n^{\text{M2L}}(\varphi) = \mathcal{L}_n(\text{mkRE}_n^{\text{M2L}}(\varphi) \cap \text{WF}_n(\varphi)) \setminus \{\emptyset\}$, and
- $\mathcal{L}_n^{\text{WS1S}}(\varphi) = \mathcal{L}_n(\text{mkRE}_n^{\text{WS1S}}(\varphi) \cap \text{WF}_n(\varphi))$.

The proof is by structural induction on φ . Above, we have seen the argument for the base case $Q a m$, other base cases follow similarly. The cases $\exists \varphi$ and $\exists \exists \varphi$ follow easily from the semantics of Π given by our concrete instantiation for π and Σ_n and the induction hypothesis. The most interesting cases are, somehow unexpectedly, those for Boolean operators. Although the definitions are purely structural, sets of encodings of models must be composed or, even worse, complemented in the inductive steps. The key property required here is that enc^{M2L} (and enc^{WS1S}) do not identify models and non-models: two different well-formed interpretations for a formula—one being a model, the other being a non-model—are encoded into different words (sets of words). This is again established by structural induction on formulas for both semantics.

Lemma 12

Let (w_1, \mathcal{I}_1) and (w_2, \mathcal{I}_2) be two M2L-well-formed interpretations for a formula φ such that $\text{enc}^{\text{M2L}}(w_1, \mathcal{I}_1) = \text{enc}^{\text{M2L}}(w_2, \mathcal{I}_2)$. Then $(w_1, \mathcal{I}_1) \models \varphi \leftrightarrow (w_2, \mathcal{I}_2) \models \varphi$.

Let (w_1, \mathcal{I}_1) and (w_2, \mathcal{I}_2) be two WS1S-well-formed interpretations for a formula φ such that $\text{enc}^{\text{WS1S}}(w_1, \mathcal{I}_1) = \text{enc}^{\text{WS1S}}(w_2, \mathcal{I}_2)$. Then $(w_1, \mathcal{I}_1) \models^{\text{B}} \varphi \leftrightarrow (w_2, \mathcal{I}_2) \models^{\text{B}} \varphi$.

5.5 Deciding language equivalence of formulas

The algorithms $\text{eqv}^{\text{M2L}} :: \mathbb{N} \rightarrow \alpha \Phi \rightarrow \alpha \Phi \rightarrow \mathbb{B}$ and $\text{eqv}^{\text{WS1S}} :: \mathbb{N} \rightarrow \alpha \Phi \rightarrow \alpha \Phi \rightarrow \mathbb{B}$ that decide language equivalence of MSO formulas check wellformedness of the

input formulas, translate the formulas into regular expressions and let eqv^{RE} do the work:

$$\text{eqv}_n^{\text{M2L}} \varphi \psi = \text{wf}_n^\Phi(\varphi \vee \psi) \wedge \text{eqv}_n^{\text{RE}}(\text{mkRE}_n^{\text{M2L}}(\varphi) + \mathbf{1})(\text{mkRE}_n^{\text{M2L}}(\psi) + \mathbf{1})$$

$$\text{eqv}_n^{\text{WS1S}} \varphi \psi = \text{wf}_n^\Phi(\varphi \vee \psi) \wedge \text{eqv}_n^{\text{RE}}(\text{mkRE}_n^{\text{WS1S}}(\varphi))(\text{mkRE}_n^{\text{WS1S}}(\psi))$$

Note that wellformedness is checked on the disjunction of both formulas to ensure that they agree on free variables (i.e., no first-order free variable of φ is used as a second-order free variable in ψ and vice versa). Further, we add the empty word into both regular expression when working with the M2L semantics. This is allowed, since \square is not a valid encoding of an interpretation, and necessary because Theorem 11 does not give us any information whether the empty word is contained in the output of mkRE^{M2L} or not.

Termination of eqv^{RE} is ensured by Theorem 6 and the definition principle of primitive recursion for wf^Φ , mkRE^{M2L} , and $\text{mkRE}^{\text{WS1S}}$. Soundness and completeness follow easily from Theorems 7, 8, and 11.

Theorem 13 (Soundness)

Let φ and ψ be MSO formulas.

- If $\text{eqv}_n^{\text{M2L}} \varphi \psi$, then $\mathcal{L}_n^{\text{M2L}}(\varphi) = \mathcal{L}_n^{\text{M2L}}(\psi)$.
- If $\text{eqv}_n^{\text{WS1S}} \varphi \psi$, then $\mathcal{L}_n^{\text{WS1S}}(\varphi) = \mathcal{L}_n^{\text{WS1S}}(\psi)$.

Theorem 14 (Completeness)

Let $\varphi \vee \psi$ be an n -well-formed MSO formula.

- If $\mathcal{L}_n^{\text{M2L}}(\varphi) = \mathcal{L}_n^{\text{M2L}}(\psi)$, then $\text{eqv}_n^{\text{M2L}} \varphi \psi$.
- If $\mathcal{L}_n^{\text{WS1S}}(\varphi) = \mathcal{L}_n^{\text{WS1S}}(\psi)$, then $\text{eqv}_n^{\text{WS1S}} \varphi \psi$.

As a sanity check let us apply our translation for M2L to the formula $\varphi = \exists (\forall 0 \in 1)$ (with names: $\exists X. \forall x. x \in X$), that is valid under the M2L semantics (but unsatisfiable under the WS1S semantics as discussed earlier). Since φ is closed, it is 0-wellformed and our underlying alphabet is $\Sigma_0 = \Sigma \times \mathbb{B}^0$ for some base alphabet Σ . For example, we can take $\Sigma = \{a\}$ and write the unique element of Σ_0 as \hat{a} . The function $\llbracket \text{mkRE}_0^{\text{M2L}}(\varphi) \rrbracket$ translates φ to the accepting Π -extended regular expression Πr over Σ_0 where r is an abbreviation:

$$r = \neg \Pi (((\text{Anth } 0 \perp)^* \cdot \text{Anth } 0 \top \cdot (\text{Anth } 0 \perp)^*) \cap \neg (\neg \mathbf{0} \cdot \text{Anth}_2 0 1 \cdot \neg \mathbf{0}))$$

Derivatives of Πr by words of the form \hat{a}^n for $n > 0$ are all \approx -equivalent to a single (also accepting) expression. More precisely, for all $w \in \Sigma_0^* \setminus \{\square\}$, we have

$$\mathcal{D}_w^*(\Pi r) = \Pi r + \Pi (r \cap \neg \Pi (((\text{Anth } 0 \perp)^*) \cap \neg (\neg \mathbf{0} \cdot \text{Anth}_2 0 1 \cdot \neg \mathbf{0})))$$

Because all derivatives of its translation are accepting, the formula φ must be valid. We would have loved to include the same example using the WS1S semantics as well, but unfortunately the output of the translation (and normalization) is a regular expression with more than 2,000 constructors (which the decision procedure still can handle).

6 Application: finite-word LTL

We want to execute the code generated by Isabelle/HOL for our decision procedures on some larger examples. For simplicity, we first focus on M2L.

In order to create larger formulas, it is helpful to introduce some syntactic abbreviations. We define the unsatisfiable formula \perp as $\exists 0 < 0$ and the valid formula \top as $\neg \perp$. Now, checking that a formula is valid amounts to checking its equivalence to \top . We call the function that performs this check *Thm*. Implication $\varphi \rightarrow \psi$ is defined as $(\neg \varphi) \vee \psi$ and universal quantification $\forall \varphi$ as before as $\neg \exists \neg \varphi$. Next, we introduce temporal logical operators *always* $\Box P :: \mathbb{N} \rightarrow \alpha \Phi$ and *eventually* $\Diamond P :: \mathbb{N} \rightarrow \alpha \Phi$ depending on $P :: \mathbb{N} \rightarrow \alpha \Phi$ —a formula parameterized by a single variable indicating the time. The operators have their usual meaning except that with the given M2L semantics the time variable ranges over a fixed set determined by the interpretation. Additionally, we lift the disjunction and implication to time-parameterized formulas.

$$\begin{aligned} \Box P t &= \forall (\neg t + 1 < 0 \rightarrow P 0) \\ \Diamond P t &= \exists (\neg t + 1 < 0 \wedge P 0) \\ (P \Rightarrow Q) t &= P t \rightarrow Q t \\ (P \vee Q) t &= P t \vee Q t \end{aligned}$$

Note that $t + 1$ has nothing to do with the next time step. It is just the lifting of the de Bruijn index under a single quantifier.

Formulas of linear temporal logic also contain atomic predicates for which the interpretation must specify at which points in time they are true. This information can be encoded in two ways, which we compare in the following.

The first possibility is to encode atomic predicates in the word of the interpretation. This is done by identifying Σ with the powerset \mathcal{P} of atomic predicates. For every point in time, that is for every position in the word, the letter is the set of predicates that are true at this point. Using this encoding, we can prove the validity of the following closed formulas over the alphabet $\mathcal{P}\{P\} = \{\{P\}, \{\}\}$ automatically within a few milliseconds.

$$\begin{aligned} \forall (\Box(Q\{P\}) \Rightarrow \Diamond(Q\{P\})) 0 \\ \forall (\Box(Q\{P\}) \Rightarrow \Box\Diamond(Q\{P\})) 0 \end{aligned}$$

Alternatively, a free second-order variable can be used to encode an atomic predicate directly. The variable denotes the set of points in time for which the atomic predicate holds. The alphabet Σ can then be trivial, i.e., $\Sigma = \{a\}$ for an arbitrary a . Using this encoding, the above two formulas correspond to

$$\begin{aligned} \forall (\Box(\lambda t. t \in 2) \Rightarrow \Diamond(\lambda t. t \in 2)) 0 \\ \forall (\Box(\lambda t. t \in 2) \Rightarrow \Box\Diamond(\lambda t. t \in 3)) 0 \end{aligned}$$

Both formulas have one free second-order variable 0 that is lifted when passing two or three quantifiers. The generated algorithm shows the equivalence to \top again within milliseconds.

In order to explore the limits of our decision procedure, formulas over more atomic predicates are required. Therefore, we consider the distributivity theorems of \Box over implication for both representations of atomic predicates as shown in Figure 2.

$$\begin{aligned}
\varphi_1 &= \forall (\Box(Q\{P\}) \Rightarrow \Box(Q\{P\})) 0 \\
\varphi_2 &= \forall (\Box(Q\{P_1\} \vee Q\{P_1, P_2\} \Rightarrow Q\{P_2\} \vee Q\{P_1, P_2\}) \Rightarrow \\
&\quad \Box(Q\{P_1\} \vee Q\{P_1, P_2\}) \Rightarrow \Box(Q\{P_2\} \vee Q\{P_1, P_2\})) 0 \\
\varphi_3 &= \forall (\Box(Q\{P_1\} \vee Q\{P_1, P_2\} \vee Q\{P_1, P_3\} \vee Q\{P_1, P_2, P_3\} \Rightarrow \\
&\quad Q\{P_2\} \vee Q\{P_1, P_2\} \vee Q\{P_2, P_3\} \vee Q\{P_1, P_2, P_3\} \Rightarrow \\
&\quad Q\{P_3\} \vee Q\{P_1, P_3\} \vee Q\{P_2, P_3\} \vee Q\{P_1, P_2, P_3\}) \Rightarrow \\
&\quad \Box(Q\{P_1\} \vee Q\{P_1, P_2\} \vee Q\{P_1, P_3\} \vee Q\{P_1, P_2, P_3\}) \Rightarrow \\
&\quad \Box(Q\{P_2\} \vee Q\{P_1, P_2\} \vee Q\{P_2, P_3\} \vee Q\{P_1, P_2, P_3\}) \Rightarrow \\
&\quad \Box(Q\{P_3\} \vee Q\{P_1, P_3\} \vee Q\{P_2, P_3\} \vee Q\{P_1, P_2, P_3\})) 0 \\
\psi_1 &= \forall (\Box(\lambda t. t \in 2) \Rightarrow \Box(\lambda t. t \in 2)) 0 \\
\psi_2 &= \forall (\Box(\lambda t. t \in 2 \rightarrow t \in 3) \Rightarrow \Box(\lambda t. t \in 2) \Rightarrow \Box(\lambda t. t \in 3)) 0 \\
\psi_3 &= \forall (\Box(\lambda t. t \in 2 \rightarrow t \in 3 \rightarrow t \in 4) \Rightarrow \Box(\lambda t. t \in 2) \Rightarrow \Box(\lambda t. t \in 3) \Rightarrow \Box(\lambda t. t \in 4)) 0
\end{aligned}$$

Fig. 2. Definition of φ_n and ψ_n .

n	size	ICFP 2013	Thm	Thm _{interm}	Thm _{dual}
1	54/ 31907	0/–	0/230.9	0/ 17.1	0/ 0.4
2	100/ 60045	2/–	0.7/ –	0/711.6	0/21.1
3	192/114857	4860/–	964.4/ –	30.4/ –	31.7/ –

Fig. 3. Benchmarks for φ_n (under M2L/WS1S semantics).

When the number of predicates n is increased, the size of φ_n grows exponentially: to express that a predicate P holds at some position we need the disjunction of all atoms containing P . In contrast, the size of ψ_n grows linearly. The complexity of ψ_n is hidden in the number of variables and therefore in its encoding—the latter also grows exponentially with increasing n .

Both, φ_i and ψ_i are theorems under both semantics. The running times of the decision procedure Thm in seconds are summarized in Figures 3 and 4 (column Thm, first number refers to the M2L semantics, the second to WS1S). Thereby, ψ_1 , ψ_2 , and ψ_3 were processed over $\Sigma = \{a\}$, φ_1 was processed over $\Sigma = \mathcal{P}\{P\}$, φ_2 over $\Sigma = \mathcal{P}\{P_1, P_2\}$ and finally φ_3 over $\Sigma = \mathcal{P}\{P_1, P_2, P_3\}$. The column “ICFP 2013” recapitulates the running times from the earlier unoptimized version of this procedure (Traytel & Nipkow, 2013).

Figures 3 and 4 also shows the sizes (column size counting the number of constructors) of the regular expressions generated from the input formulas. These numbers show a huge gap between WS1S and M2L that also shows up in the runtime results. Our implementation of \mathcal{Q} is very inefficient. As future work, we plan

n	size	ICFP 2013	Thm	Thm _{interm}	Thm _{dual}
1	54/31907	0/–	0/224.4	0/ 22.6	0/ 0.4
2	81/48797	2/–	0.3/ –	0/434.1	0/14.4
3	108/65687	2640/–	64.4/ –	3.9/ –	17.8/ –

Fig. 4. Benchmarks for ψ_n (under M2L/WS1S semantics).

to investigate the addition of this regular operator as a constructor to the data type of regular expressions, similarly to our addition of the projection operator.

The last two columns show the running times of two variations of Thm : $\text{Thm}_{\text{interm}}$ and Thm_{dual} . One remaining source of inefficiency in Thm is the fact that, although it constructs a bisimulation modulo $\langle\!\langle - \rangle\!\rangle$, the intermediate expressions on which the derivatives are computed are only ACI-normalized (i.e., they might contain redundant subexpressions like $\mathbf{0} \cap r$ and the derivative would needlessly recurse in r). The algorithm $\text{Thm}_{\text{interm}}$ addresses this inefficiency by normalizing the intermediate expressions with $\langle\!\langle - \rangle\!\rangle$. This intermediate normalization might seem harmless, but it is not clear anymore that the number of derivatives interspersed with normalization is finite³. We were not able to prove finiteness of such derivatives interspersed with $\langle\!\langle - \rangle\!\rangle$ (although we conjecture that it holds). However, we have proved that under the condition that $\text{Thm}_{\text{interm}}$ terminates, its output—namely \top if the input formula is valid, \perp otherwise—is correct.

The algorithm Thm_{dual} is similar to $\text{Thm}_{\text{interm}}$ in the respect that it normalizes intermediate expressions. Therefore, we again guarantee only partial correctness. Unlike $\text{Thm}_{\text{interm}}$, Thm_{dual} works with dual regular expressions (Section 4.5). It seems to be the better choice for the WS1S semantics.

The attentive reader will have noticed that we have said nothing about how sets are represented in the code generated from our mathematical definitions. We use the default implementation as lists (with a linear membership test) from Isabelle’s library for our measurements. We have also experimented with an existing verified red–black tree implementation. Isabelle’s code generator supports the transparent replacement of sets by some verified implementation (Haftmann *et al.*, 2013). Unfortunately, the overhead incurred by the trees outweighed the gain of a logarithmic membership test instead of a linear one.

The performance of our automatically generated code may appear disappointing but that would be a misunderstanding of our intentions. We see our work primarily as a succinct and elegant functional program that may pave the way towards verified and efficient decision procedures. As a bonus, the generated code is applicable to small examples. In the context of interactive theorem proving, this is primarily what one encounters: small formulas. Any automation is welcome here because it saves the user time and effort. Automatic verification of larger systems is the domain of highly tuned implementations such as MONA.

7 Conclusion

We have presented functional programs that decide equivalence of MSO formulas for two different semantics in Isabelle/HOL. They come with formal proofs of termination, soundness and completeness. The programs operate by translating formulas into Π -extended regular expressions and deciding the language equivalence

³ For example, using the terminating normalization function that does the same ACI simplifications as $\langle\!\langle - \rangle\!\rangle$, but additionally soundly rewrites $\mathbf{1} \cdot a^*$ to $a^* \cdot a^*$ for a fixed symbol a will result in an infinite number of derivatives when applied at intermediate steps to the initial expression a^* .

of the latter using Brzozowski derivatives. Although formalized in Isabelle/HOL's functional programming language, we can automatically generate code from them in different functional target languages. The development amounts to roughly 500 lines of functional programs and 6,500 lines of proofs, of which 3,000 lines are devoted to deciding equivalence of Π -extended regular expressions. The functional programs are completely contained in this paper. The Isabelle scripts are publicly available (Traytel & Nipkow, 2014).

Our work can be continued in two dimensions. First, our algorithm still offers much room for optimization. Especially, the inefficient formalization of \mathcal{Q} should be revised. Second, several related decidable logics can be formalized and verified using similar technology. A related logic is MSO on infinite words (also called S1S). S1S formulas can be translated into ω -regular expressions representing ω -regular languages. A verified decision procedure for deciding equivalence of ω -regular expressions without constructing ω -automata is an interesting challenge. A similarly ambitious goal is to move from words to trees (or even from ω -words to ω -trees) and decide equivalence of MSO formulas on (in)finite trees (or alternatively (W)S2S formulas) by translating them into (ω -)regular tree expressions.

Acknowledgments

We thank Alexander Krauss for inspiring discussions, Jasmin Blanchette for numerous comments on the presentation, and several anonymous reviewers for a wealth of comments and questions that helped to clarify certain fine points. While carrying out this work, Traytel was affiliated with TU München and supported by the doctorate program 1480 (PUMA) of the Deutsche Forschungsgemeinschaft (DFG).

References

- Antimirov, V. (1996) Partial derivatives of regular expressions and finite automaton constructions. *Theor. Comput. Sci.* **155**(2), 291–319.
- Asperti, A. (2012) A compact proof of decidability for regular expression equivalence. In *Proc. Int. Conf. Interactive Theorem Proving, ITP 2012*, Beringer, L. & Felty, A. (eds), Lect. Notes Comput. Sci., vol. 7406. Springer, pp. 283–298.
- Ayari, A. & Basin, D. (2000) Bounded model construction for monadic second-order logics. In *Proc. Int. Conf. Computer Aided Verification, CAV 2000*, Emerson, E. A. & Sistla, A. P. (eds), Lect. Notes Comput. Sci., vol. 1855. Springer, pp. 99–112.
- Ballarin, C. (2006) Interpretation of locales in Isabelle: Theories and proof contexts. In *Proc. Int. Conf. Mathematical Knowledge Management, MKM 2006*, Borwein, J. M. & Farmer, W. M. (eds), Lect. Notes Comput. Sci., vol. 4108. Springer, pp. 31–43.
- Basin, D. & Friedrich, S. (2000) Combining WS1S and HOL. In *Frontiers of Combining Systems 2*, Gabbay, D. M. & de Rijke, M. (eds), Studies in Logic and Computation, vol. 7. Research Studies Press, pp. 39–56.
- Berghofer, S. & Reiter, M. (2009) Formalizing the logic-automaton connection. In *Proc. Int. Conf. Theorem Proving in Higher Order Logics, TPHOLs 2009*, Berghofer, S., Nipkow, T., Urban, C. & Wenzel, M. (eds), Lect. Notes Comput. Sci., vol. 5674. Springer, pp. 147–163.

- Bonchi, F. & Pous, D. (2013) Checking NFA equivalence with bisimulations up to congruence. In *Proc. Int. Symp. Principles of Programming Languages, POPL 2013*, Giacobazzi, R. & Cousot, R. (eds), ACM, pp. 457–468.
- Braibant, T. & Pous, D. (2010) An efficient Coq tactic for deciding Kleene algebras. In *Proc. Int. Conf. Interactive Theorem Proving, ITP 2010*, Kaufmann, M. & Paulson, L. (eds), Lect. Notes Comput. Sci., vol. 6172. Springer, pp. 163–178.
- Brzozowski, J. A. (1964) Derivatives of regular expressions. *J. ACM* **11**(4), 481–494.
- Caron, P., Champarnaud, J.-M., & Mignot, L. (2011) Partial derivatives of an extended regular expression. In *Proc. Int. Conf. Language and Automata Theory and Applications, LATA 2011*, Dediu, A.-H., Inenaga, S. & Martín-Vide, C. (eds), Lect. Notes Comput. Sci., vol. 6638. Springer, pp. 179–191.
- Coquand, T. & Siles, V. (2011) A decision procedure for regular expression equivalence in type theory. In *Proc. Int. Conf. Certified Programs and Proofs, CPP 2011*, Jouannaud, J.-P. & Shao, Z. (eds), Lect. Notes Comput. Sci., vol. 7086. Springer, pp. 119–134.
- Danielsson, N. A. (2010) Total parser combinators. In *Proc. Int. Conf. Functional Programming, ICFP 2010*, Hudak, P. & Weirich, S. (eds), ACM, pp. 285–296.
- Elgaard, J., Klarlund, N. & Møller, A. (1998) MONA 1.x: New techniques for WS1S and WS2S. In *Proc. Int. Conf. Computer Aided Verification, CAV 1998*, Hu, A. J. & Vardi, M. Y. (eds), Lect. Notes Comput. Sci., vol. 1427. Springer, pp. 516–520.
- Fischer, S., Huch, F. & Wilke, T. (2010) A play on regular expressions: Functional pearl. *Proc. Int. Conf. Functional Programming, ICFP 2010*, Hudak, P. & Weirich, S. (eds), ACM, pp. 357–368.
- Ginzburg, A. (1967) A procedure for checking equality of regular expressions. *J. ACM* **14**(2), 355–362.
- Haftmann, F. & Nipkow, T. (2010) Code generation via higher-order rewrite systems. *Proc. Int. Symp. Functional and Logic Programming, FLOPS 2010*, Lect. Notes Comput. Sci., vol. 6009. Springer, pp. 103–117.
- Haftmann, F., Krauss, A., Kunčar, O. & Nipkow, T. (2013) Data refinement in Isabelle/HOL. In *Proc. Int. Conf. Interactive Theorem Proving, ITP 2013*, Blazy, S., Paulin-Mohring, C. & Pichardie, D. (eds), Lect. Notes Comput. Sci., vol. 7998. Springer, pp. 100–115.
- Henriksen, J. G., Jensen, J. L., Jørgensen, M. E., Klarlund, N., Paige, R., Rauhe, T. & Sandholm, A. (1995) MONA: Monadic second-order logic in practice. In *Proc. Int. Workshop Tools and Algorithms for the Construction and Analysis of Systems, TACAS 1995*, Brinksma, E., Cleaveland, R., Larsen, K., Margaria, T. & Steffen, B. (eds), Lect. Notes Comput. Sci., vol. 1019. Springer, pp. 89–110.
- Klarlund, N. (1999) A theory of restrictions for logics and automata. In *Proc. Int. Conf. Computer Aided Verification, CAV 1999*, Halbwachs, N. & Peled, D. (eds), Lect. Notes Comput. Sci., vol. 1633. Springer, pp. 406–417.
- Kozen, D. (2008 March) *On the Coalgebraic Theory of Kleene Algebra with Tests*. Tech. rept. <http://hdl.handle.net/1813/10173>. Computing and Information Science, Cornell University.
- Krauss, A. & Nipkow, T. (2012) Proof pearl: Regular expression equivalence and relation algebra. *J. Autom. Reason.* **49**(1), 95–106. published online March 2011.
- Might, M., Darais, D. & Spiewak, D. (2011) Parsing with derivatives: A functional pearl. In *Proc. Int. Conf. Functional Programming, ICFP 2011*, Chakravarty, M. M. T., Hu, Z. & Danvy, O. (eds), ACM, pp. 189–195.
- Moreira, N., Pereira, D. & de Sousa, S. M. (2012) Deciding regular expressions (in-)equivalence in Coq. In *Relational and Algebraic Methods in Computer Science, RAMiCS 2012*, Kahl, W. & Griffin, T. (eds), Lect. Notes Comput. Sci., vol. 7560. Springer, pp. 98–113.

- Nipkow, T. & Klein, G. (2014) *Concrete Semantics: With Isabelle/HOL*. Springer. Available at: <http://www.in.tum.de/~nipkow/Concrete-Semantics>.
- Nipkow, T. & Traytel, D. (2014) Unified decision procedures for regular expression equivalence. In *Proc. Int. Conf. Interactive Theorem Proving, ITP 2014*, Klein, G. & Gamboa, R. (eds), Lect. Notes Comput. Sci., vol. 8558. Springer, pp. 450–466.
- Nipkow, T., Paulson, L. & Wenzel, M. (2002) *Isabelle/HOL — A Proof Assistant for Higher-Order Logic*. Lect. Notes Comput. Sci., vol. 2283. Springer.
- Okhotin, A. (2005) The dual of concatenation. *Theor. Comput. Sci.* **345**(2–3), 425–447.
- Owens, S., Reppy, J. H. & Turon, A. (2009) Regular-expression derivatives re-examined. *J. Funct. Program.* **19**(2), 173–190.
- Owre, S. & Rueß, H. (2000) Integrating WS1S with PVS. In *Proc. Int. Conf. Computer Aided Verification, CAV 2000*, Emerson, E. A. & Sistla, A. P. (eds), Lect. Notes Comput. Sci., vol. 1855. Springer, pp. 548–551.
- Pous, D. (2015) Symbolic algorithms for language equivalence and Kleene algebra with test. In *Proc. Int. Symp. Principles of Programming Languages, POPL 2015*, Walker, D. (ed), ACM, pp. 357–368.
- Rutten, Jan J. M. M. (1998) Automata and coinduction (an exercise in coalgebra) In *Proc. Int. Conf. Concurrency Theory, CONCUR 1998*, Sangiorgi, D. & de Simone, R. (eds), Lect. Notes Comput. Sci., vol. 1466. Springer, pp. 194–218.
- Thomas, W. (1997) Languages, automata, and logic. In *Handbook of Formal Languages*, Rozenberg, G. & Salomaa, A. (eds), Springer, pp. 389–455.
- Traytel, D. & Nipkow, T. (2013) Verified decision procedures for MSO on words based on derivatives of regular expressions. *Proc. Int. Conf. Functional Programming, ICFP 2013*, Morrisett, G. & Uustalu, T. (eds), ACM, pp. 3–12.
- Traytel, D. & Nipkow, T. (2014) Decision procedures for MSO on words based on derivatives of regular expressions. In *Archive of Formal Proofs*, Klein, G., Nipkow, T. & Paulson, L. (eds), http://afp.sf.net/entries/MSO_Regex_Equivalence.shtml, Formal proof development.
- Wu, C., Zhang, X. & Urban, C. (2014) A formalisation of the Myhill-Nerode theorem based on regular expressions. *J. Autom. Reason.* **52**(4), 451–480.