

# A Formalisation of the Myhill-Nerode Theorem based on Regular Expressions<sup>\*</sup>

Chunhan Wu

PLA University of Science and Technology Nanjing, China  
and

Xingyuan Zhang

PLA University of Science and Technology Nanjing, China  
and

Christian Urban

King's College London, United Kingdom

---

There are numerous textbooks on regular languages. Many of them focus on finite automata for proving properties. Unfortunately, automata are not so straightforward to formalise in theorem provers. The reason is that natural representations for automata are graphs, matrices or functions, none of which are inductive datatypes. Regular expressions can be defined straightforwardly as a datatype and a corresponding reasoning infrastructure comes for free in theorem provers. We show in this paper that a central result from formal language theory—the Myhill-Nerode Theorem—can be recreated using only regular expressions. From this theorem many closure properties of regular languages follow.

Categories and Subject Descriptors: F.4.1 [**Mathematical Logic and Formal Languages**]: Mechanical Theorem Proving; F.4.3 [**Mathematical Logic and Formal Languages**]: Formal Languages

General Terms: Interactive theorem proving, regular languages

Additional Key Words and Phrases: Myhill-Nerode theorem, regular expressions, Isabelle theorem prover

---

## 1. INTRODUCTION

Regular languages are an important and well-understood subject in Computer Science, with many beautiful theorems and many useful algorithms. There is a wide range of textbooks on this subject, many of which are aimed at students and contain very detailed ‘pencil-and-paper’ proofs (e.g. the textbooks by Hopcroft and Ullman [1969] and by Kozen [1997]). It seems natural to exercise theorem provers by formalising the theorems and by verifying formally the algorithms.

A popular choice for a theorem prover would be one based on Higher-Order Logic

---

Corresponding Author: Christian Urban, Department of Informatics, King's College London, Strand, London WC2R 2LS, UK. Email: christian.urban@kcl.ac.uk.

<sup>\*</sup>This is a revised and expanded version of [Wu et al. 2011a].

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 20YY ACM 1529-3785/YY/00-0001 \$5.00

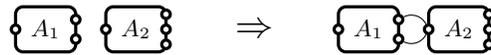
(HOL), for example HOL4, HOLlight or Isabelle/HOL. For the development presented in this paper we will use the Isabelle/HOL. HOL is a predicate calculus that allows quantification over predicate variables. Its type system is based on the Simple Theory of Types by Church [1940]. Although many mathematical concepts can be conveniently expressed in HOL, there are some limitations that hurt when attempting a simple-minded formalisation of regular languages in it.

The typical approach to regular languages, taken for example by Hopcroft and Ullman [1969] and by Kozen [1997], is to introduce finite deterministic automata and then define most notions in terms of them. For example, a regular language is normally defined as:

**DEFINITION 1.1.** *A language  $A$  is regular, provided there is a finite deterministic automaton that recognises all strings of  $A$ .*

This approach has many benefits. Among them is the fact that it is easy to convince oneself that regular languages are closed under complementation: one just has to exchange the accepting and non-accepting states in the corresponding automaton to obtain an automaton for the complement language. The problem, however, lies with formalising such reasoning in a theorem prover. Automata are built up from states and transitions that are commonly represented as graphs, matrices or functions, none of which, unfortunately, can be defined as an inductive datatype.

In case of graphs and matrices, this means we have to build our own reasoning infrastructure for them, as neither Isabelle/HOL nor HOL4 nor HOLlight support them with libraries. Also, reasoning about graphs and matrices can be a hassle in theorem provers, because we have to be able to combine automata. Consider for example the operation of sequencing two automata, say  $A_1$  and  $A_2$ , by connecting the accepting states of  $A_1$  to the initial state of  $A_2$ :



On ‘paper’ we can define the corresponding graph in terms of the disjoint union of the state nodes. Unfortunately in HOL, the standard definition for disjoint union, namely

$$A_1 \uplus A_2 \stackrel{\text{def}}{=} \{(1, x) \mid x \in A_1\} \cup \{(2, y) \mid y \in A_2\} \quad (1)$$

changes the type—the disjoint union is not a set, but a set of pairs. Using this definition for disjoint union means we do not have a single type for the states of automata. As a result we will not be able to define a regular language as one for which there exists an automaton that recognises all its strings (Definition 1.1). This is because we cannot make a definition in HOL that is only polymorphic in the state type, but not in the predicate for regularity; and there is no type quantification available in HOL (unlike in Coq, for example).<sup>1</sup>

An alternative, which provides us with a single type for states of automata, is to give every state node an identity, for example a natural number, and then be careful to rename these identities apart whenever connecting two automata. This results in clunky proofs establishing that properties are invariant under renaming. Similarly, connecting two automata represented as matrices results in messy constructions, which are not pleasant to formally reason about. Braibant [2012, Page 67], for example, writes that there are no

<sup>1</sup>Slind already pointed out this problem in an email to the HOL4 mailing list on 21st April 2005.

problems with reasoning about matrices, but that there is an “intrinsic difficulty of working with rectangular matrices” in some parts of his formalisation of formal languages in Coq.

Functions are much better supported in Isabelle/HOL, but they still lead to similar problems as with graphs. Composing, for example, two non-deterministic automata in parallel requires also the formalisation of disjoint unions. Nipkow [1998] dismisses for this the option of using identities, because it leads according to him to “messy proofs”. Since he does not need to define what regular languages are, Nipkow opts for a variant of (1) using bit lists, but writes

*“All lemmas appear obvious given a picture of the composition of automata. . . Yet their proofs require a painful amount of detail.”*

and

*“If the reader finds the above treatment in terms of bit lists revoltingly concrete, I cannot disagree. A more abstract approach is clearly desirable.”*

Because of these problems to do with representing automata, formalising automata theory is surprisingly not as easy as one might think, despite the sometimes very detailed, but informal, textbook proofs. Lammich and Tuerk [2012] formalised Hopcroft’s algorithm using an automata library of 27 kloc in Isabelle/HOL. There they use matrices for representing automata. Functions are used by Nipkow [1998], who establishes the link between regular expressions and automata in the context of lexing. Berghofer and Reiter [2009] use functions as well for formalising automata working over bit strings in the context of Presburger arithmetic. A Larger formalisation of automata theory, including the Myhill-Nerode theorem, was carried out in Nuprl by Constable et al. [2000] which also uses functions. Other large formalisations of automata theory in Coq are by Filliâtre [1997] who essentially uses graphs and by Almeida et al. [2010] and Braibant [2012], who both use matrices. Many of these works, like Nipkow [1998] or Braibant [2012], mention intrinsic problems with their representation of automata which made them ‘fight’ their respective theorem prover.

In this paper, we will not attempt to formalise automata theory in Isabelle/HOL nor will we attempt to formalise automata proofs from the literature, but take a different approach to regular languages than is usually taken. Instead of defining a regular language as one where there exists an automaton that recognises all its strings, we define a regular language as:

**DEFINITION 1.2.** *A language  $A$  is regular, provided there is a regular expression that matches all strings of  $A$ .*

And then ‘forget’ automata completely. The reason is that regular expressions can be defined as an inductive datatype and a reasoning infrastructure for them (like induction and recursion) comes for free in HOL. This convenience of regular expressions has recently been exploited in HOL4 with a formalisation of regular expression matching based on derivatives by Owens and Slind [2008], and with an equivalence checker for regular expressions in Isabelle/HOL by Krauss and Nipkow [2012] and in Matida by Asperti [2012] and in Coq by Coquand and Siles [2011]. The main purpose of this paper is to show that a central result about regular languages—the Myhill-Nerode Theorem—can be recreated by only using regular expressions. This theorem gives necessary and sufficient conditions for

when a language is regular. As a corollary of this theorem we can easily establish the usual closure properties, including complementation, for regular languages. We use the Continuation Lemma, which is also a corollary of the Myhill-Nerode Theorem, for establishing the non-regularity of the language  $a^n b^n$ .

**Contributions:** There is an extensive literature on regular languages. To our best knowledge, our proof of the Myhill-Nerode Theorem is the first that is based on regular expressions, only. The part of this theorem stating that finitely many partitions imply regularity of the language is proved by an argument about solving equational systems. This argument appears to be folklore. For the other part, we give two proofs: one direct proof using certain tagging-functions, and another indirect proof using Antimirov’s partial derivatives [1995]. Again to our best knowledge, the tagging-functions have not been used before for establishing the Myhill-Nerode Theorem. Derivatives of regular expressions have been used recently quite widely in the literature; partial derivatives, in contrast, attract much less attention. However, partial derivatives are more suitable in the context of the Myhill-Nerode Theorem, since it is easier to establish formally their finiteness result. We are not aware of any proof that uses either of them for proving the Myhill-Nerode Theorem.

## 2. PRELIMINARIES

Strings in Isabelle/HOL are lists of characters with the *empty string* being represented by the empty list, written  $[]$ . We assume there are only finitely many different characters. *Languages* are sets of strings. The language containing all strings is written in Isabelle/HOL as *UNIV*. The concatenation of two languages is written  $A \cdot B$  and a language raised to the power  $n$  is written  $A^n$ . They are defined as usual

$$\begin{aligned} A \cdot B &\stackrel{\text{def}}{=} \{s_1 @ s_2 \mid s_1 \in A \wedge s_2 \in B\} \\ A^0 &\stackrel{\text{def}}{=} \{[]\} \\ A^{n+1} &\stackrel{\text{def}}{=} A \cdot A^n \end{aligned}$$

where  $@$  is the list-append operation. The Kleene-star of a language  $A$  is defined as the union over all powers, namely  $A^* \stackrel{\text{def}}{=} \bigcup_n A^n$ . In the paper we will make use of the following properties of these constructions.

**PROPOSITION 2.1.**

- (i)  $A^* = A \cdot A^* \cup \{[]\}$
- (ii) If  $[] \notin A$  and  $s \in A^{n+1}$  then  $n < \text{length } s$ .
- (iii)  $B \cdot (\bigcup_n A^n) = (\bigcup_n B \cdot A^n)$
- (iv) If  $x \in A^*$  and  $x \neq []$  then there exists an  $x_p$  and  $x_s$  with  $x = x_p @ x_s$  and  $x_p \neq []$  such that  $x_p \in A$  and  $x_s \in A^*$ .

In (ii) we use the notation *length*  $s$  for the length of a string; this property states that if  $[] \notin A$  then the lengths of the strings in  $A^{n+1}$  must be longer than  $n$ . Property (iv) states that a non-empty string in  $A^*$  can always be split up into a non-empty prefix belonging to  $A$  and the rest being in  $A^*$ . We omit the proofs for these properties, but invite the reader to consult our formalisation.<sup>2</sup>

<sup>2</sup>Available under Wu et al. [2011b] in the Archive of Formal Proofs at <http://afp.sourceforge.net/entries/Myhill-Nerode.shtml>.

The notation in Isabelle/HOL for the quotient of a language  $A$  according to an equivalence relation  $\approx$  is  $A // \approx$ . We will write  $\llbracket x \rrbracket_{\approx}$  for the equivalence class defined as  $\{y \mid y \approx x\}$ , and have  $x \approx y$  if and only if  $\llbracket x \rrbracket_{\approx} = \llbracket y \rrbracket_{\approx}$ .

Central to our proof will be the solution of equational systems involving equivalence classes of languages. For this we will use Arden's Lemma (see for example [Sakarovitch 2009, Page 100]), which solves equations of the form  $X = A \cdot X \cup B$  provided  $\square \notin A$ . However we will need the following 'reversed' version of Arden's Lemma ('reversed' in the sense of changing the order of  $A \cdot X$  to  $X \cdot A$ ).<sup>3</sup>

LEMMA 2.1 (REVERSED ARDEN'S LEMMA).

If  $\square \notin A$  then  $X = X \cdot A \cup B$  if and only if  $X = B \cdot A^*$ .

Regular expressions are defined as the inductive datatype

$$\begin{array}{l}
r ::= \text{ZERO} \\
\quad | \text{ONE} \\
\quad | \text{ATOM } c \\
\quad | \text{TIMES } r \ r \\
\quad | \text{PLUS } r \ r \\
\quad | \text{STAR } r
\end{array}$$

and the language matched by a regular expression is defined by recursion as

$$\begin{array}{l}
\mathcal{L}(\text{ZERO}) \stackrel{\text{def}}{=} \{\} \\
\mathcal{L}(\text{ONE}) \stackrel{\text{def}}{=} \{\square\} \\
\mathcal{L}(\text{ATOM } c) \stackrel{\text{def}}{=} \{[c]\} \\
\mathcal{L}(\text{PLUS } r_1 \ r_2) \stackrel{\text{def}}{=} \mathcal{L}(r_1) \cup \mathcal{L}(r_2) \\
\mathcal{L}(\text{TIMES } r_1 \ r_2) \stackrel{\text{def}}{=} \mathcal{L}(r_1) \cdot \mathcal{L}(r_2) \\
\mathcal{L}(\text{STAR } r) \stackrel{\text{def}}{=} \mathcal{L}(r)^*
\end{array}$$

Given a finite set of regular expressions  $rs$ , we will make use of the operation of generating a regular expression that matches the union of all languages of  $rs$ . This definition is not trivial in a theorem prover, because  $rs$  (being a set) is unordered, but the regular expression needs an order. Since we only need to know the existence of such a regular expression, we can use Isabelle/HOL's *fold\_graph* and Hilbert's choice operator, written *SOME* in Isabelle/HOL, for defining  $\dagger rs$ . This operation, roughly speaking, folds *PLUS* over the set  $rs$  with *ZERO* for the empty set. We can prove that for a finite set  $rs$

$$\mathcal{L}(\dagger rs) = \bigcup (\mathcal{L} \ ' \ rs) \tag{2}$$

holds, whereby  $\mathcal{L} \ ' \ rs$  stands for the image of the set  $rs$  under function  $\mathcal{L}$  defined as

$$\mathcal{L} \ ' \ rs \stackrel{\text{def}}{=} \{\mathcal{L}(r) \mid r \in rs\}$$

In what follows we shall use this convenient short-hand notation for images of sets also with other functions.

<sup>3</sup>The details of the proof for the reversed Arden's Lemma are given in the Appendix.

### 3. THE MYHILL-NERODE THEOREM, FIRST PART

The key definition in the Myhill-Nerode Theorem is the *Myhill-Nerode Relation*, which states that w.r.t. a language two strings are related, provided there is no distinguishing extension in this language. This can be defined as a ternary relation.

DEFINITION 3.1 (MYHILL-NERODE RELATION). *Given a language  $A$ , two strings  $x$  and  $y$  are Myhill-Nerode related provided*

$$x \approx_A y \stackrel{\text{def}}{=} \forall z. (x @ z \in A) = (y @ z \in A)$$

It is easy to see that  $\approx_A$  is an equivalence relation, which partitions the set of all strings,  $UNIV$ , into a set of disjoint equivalence classes. To illustrate this quotient construction, let us give a simple example: consider the regular language containing just the string  $[c]$ . The relation  $\approx_{\{[c]\}}$  partitions  $UNIV$  into three equivalence classes  $X_1, X_2$  and  $X_3$  as follows

$$\begin{aligned} X_1 &= \{\emptyset\} \\ X_2 &= \{[c]\} \\ X_3 &= UNIV - \{\emptyset, [c]\} \end{aligned}$$

One direction of the Myhill-Nerode Theorem establishes that if there are finitely many equivalence classes, like in the example above, then the language is regular. In our setting we therefore have to show:

THEOREM 3.1. *If finite ( $UNIV // \approx_A$ ) then regular  $A$ .*

To prove this theorem, we first define the set *finals*  $A$  as those equivalence classes from  $UNIV // \approx_A$  that contain strings of  $A$ , namely

$$\text{finals } A \stackrel{\text{def}}{=} \{\llbracket s \rrbracket_{\approx_A} \mid s \in A\} \quad (3)$$

In our running example,  $X_2$  is the only equivalence class in *finals*  $\{[c]\}$ . It is straightforward to show that in general

$$A = \bigcup \text{finals } A \quad \text{finals } A \subseteq UNIV // \approx_A \quad (4)$$

hold. Therefore if we know that there exists a regular expression for every equivalence class in *finals*  $A$  (which by assumption must be a finite set), then we can use  $\dagger$  to obtain a regular expression that matches every string in  $A$ .

Our proof of Theorem 3.1 relies on a method that can calculate a regular expression for every equivalence class, not just the ones in *finals*  $A$ . We first define the notion of *one-character-transition* between two equivalence classes

$$Y \stackrel{c}{\rightrightarrows} X \stackrel{\text{def}}{=} Y \cdot \{[c]\} \subseteq X \quad (5)$$

which means that if we append the character  $c$  to the end of all strings in the equivalence class  $Y$ , we obtain a subset of  $X$ . Note that we do not define an automaton here, we merely relate two sets (with the help of a character). In our concrete example we have  $X_1 \stackrel{c}{\rightrightarrows} X_2$ ,  $X_1 \stackrel{d_i}{\rightrightarrows} X_3$  with  $d_i$  being any other character than  $c$ , and  $X_3 \stackrel{c_j}{\rightrightarrows} X_3$  for any character  $c_j$ .

Next we construct an *initial equational system* that contains an equation for each equivalence class. We first give an informal description of this construction. Suppose we have the equivalence classes  $X_1, \dots, X_n$ , there must be one and only one that contains the empty string  $\emptyset$  (since equivalence classes are disjoint). Let us assume  $\emptyset \in X_1$ . We build the following initial equational system

$$\begin{aligned}
X_1 &= (Y_{11}, ATOM\ c_{11}) + \dots + (Y_{1p}, ATOM\ c_{1p}) + \lambda(ONE) \\
X_2 &= (Y_{21}, ATOM\ c_{21}) + \dots + (Y_{2o}, ATOM\ c_{2o}) \\
&\vdots \\
X_n &= (Y_{n1}, ATOM\ c_{n1}) + \dots + (Y_{nq}, ATOM\ c_{nq})
\end{aligned}$$

where the terms  $(Y_{ij}, ATOM\ c_{ij})$  are pairs consisting of an equivalence class and a regular expression. In the initial equational system, they stand for all transitions  $Y_{ij} \xrightarrow{c_{ij}} X_i$ . There can only be finitely many terms of the form  $(Y_{ij}, ATOM\ c_{ij})$  in a right-hand side since by assumption there are only finitely many equivalence classes and only finitely many characters. The term  $\lambda(ONE)$  in the first equation acts as a marker for the initial state, that is the equivalence class containing the empty string  $\epsilon$ .<sup>4</sup> In our running example we have the initial equational system

$$\begin{aligned}
X_1 &= \lambda(ONE) \\
X_2 &= (X_1, ATOM\ c) \\
X_3 &= (X_1, ATOM\ d_1) + \dots + (X_1, ATOM\ d_n) \\
&\quad + (X_3, ATOM\ c_1) + \dots + (X_3, ATOM\ c_m)
\end{aligned} \tag{6}$$

where  $d_1 \dots d_n$  is the sequence of all characters but not containing  $c$ , and  $c_1 \dots c_m$  is the sequence of all characters.

Overloading the function  $\mathcal{L}$  for the two kinds of terms in the equational system, we have

$$\mathcal{L}(Y, r) \stackrel{def}{=} Y \cdot \mathcal{L}(r) \quad \mathcal{L}(\lambda(r)) \stackrel{def}{=} \mathcal{L}(r)$$

and we can prove for  $X_2 \dots X_n$  that the following equations

$$X_i = \mathcal{L}(Y_{i1}, ATOM\ c_{i1}) \cup \dots \cup \mathcal{L}(Y_{iq}, ATOM\ c_{iq}). \tag{7}$$

hold. Similarly for  $X_1$  we can show the following equation

$$X_1 = \mathcal{L}(Y_{11}, ATOM\ c_{11}) \cup \dots \cup \mathcal{L}(Y_{1p}, ATOM\ c_{1p}) \cup \mathcal{L}(\lambda(ONE)). \tag{8}$$

holds. The reason for adding the  $\lambda$ -marker to our initial equational system is to obtain this equation: it only holds with the marker, since none of the other terms contain the empty string. The point of the initial equational system is that solving it means we will be able to extract a regular expression for every equivalence class.

Our representation for the equations in Isabelle/HOL are pairs, where the first component is an equivalence class (a set of strings) and the second component is a set of terms. Given a set of equivalence classes  $CS$ , our initial equational system  $Init\ CS$  is thus formally defined as

$$\begin{aligned}
Init\_rhs\ CS\ X &\stackrel{def}{=} \text{if } \epsilon \in X \\
&\quad \text{then } \{(Y, ATOM\ c) \mid Y \in CS \wedge Y \xrightarrow{c} X\} \cup \{\lambda(ONE)\} \\
&\quad \text{else } \{(Y, ATOM\ c) \mid Y \in CS \wedge Y \xrightarrow{c} X\} \\
Init\ CS &\stackrel{def}{=} \{(X, Init\_rhs\ CS\ X) \mid X \in CS\}
\end{aligned} \tag{9}$$

<sup>4</sup>Note that we mark, roughly speaking, the single ‘initial’ state in the equational system, which is different from the method by Brzozowski [1964], where he marks the ‘terminal’ states. We are forced to set up the equational system in our way, because the Myhill-Nerode Relation determines the ‘direction’ of the transitions—the successor ‘state’ of an equivalence class  $Y$  can be reached by adding a character to the end of  $Y$ . This is also the reason why we have to use our reversed version of Arden’s Lemma.

Because we use sets of terms for representing the right-hand sides of equations, we can prove (7) and (8) more concisely as

LEMMA 3.1. *If  $(X, rhs) \in \text{Init} (UNIV // \approx_A)$  then  $X = \bigcup \mathcal{L} \text{ ' } rhs$ .*

Our proof of Theorem 3.1 will proceed by transforming the initial equational system into one in *solved form* maintaining the invariant in Lemma 3.1. From the solved form we will be able to read off the regular expressions.

In order to transform an equational system into solved form, we have two operations: one that takes an equation of the form  $X = rhs$  and removes any recursive occurrences of  $X$  in the  $rhs$  using our variant of Arden's Lemma. The other operation takes an equation  $X = rhs$  and substitutes  $X$  throughout the rest of the equational system adjusting the remaining regular expressions appropriately. To define this adjustment we define the *append-operation* taking a term and a regular expression as argument

$$\begin{aligned} (Y, r_2) \triangleleft r_1 &\stackrel{\text{def}}{=} (Y, \text{TIMES } r_2 r_1) \\ \lambda(r_2) \triangleleft r_1 &\stackrel{\text{def}}{=} \lambda(\text{TIMES } r_2 r_1) \end{aligned}$$

We lift this operation to entire right-hand sides of equations, written as  $rhs \triangleleft r$ . With this we can define the *arden-operation* for an equation of the form  $X = rhs$  as:

$$\begin{aligned} \text{Arden } X \text{ } rhs &\stackrel{\text{def}}{=} \text{let} \\ &\quad rhs' = rhs - \{(X, r) \mid (X, r) \in rhs\} \\ &\quad r' = \text{STAR} (\text{+}\{r \mid (X, r) \in rhs\}) \\ &\text{in } rhs' \triangleleft r' \end{aligned} \quad (10)$$

In this definition, we first delete all terms of the form  $(X, r)$  from  $rhs$ ; then we calculate the combined regular expressions for all  $r$  coming from the deleted  $(X, r)$ , and take the *STAR* of it; finally we append this regular expression to  $rhs'$ . If we apply this operation to the right-hand side of  $X_3$  in (6), we obtain the equation:

$$\begin{aligned} X_3 = &(X_1, \text{TIMES } (ATOM d_1) (\text{STAR} \text{+}\{ATOM c_1, \dots, ATOM c_m\})) + \dots \\ &\dots + (X_1, \text{TIMES } (ATOM d_n) (\text{STAR} \text{+}\{ATOM c_1, \dots, ATOM c_m\})) \end{aligned}$$

That means we eliminated the recursive occurrence of  $X_3$  on the right-hand side. Note we used the abbreviation  $\text{+}\{ATOM c_1, \dots, ATOM c_m\}$  to stand for a regular expression that matches with every character. In our algorithm we are only interested in the existence of such a regular expression and do not specify it any further.

It can be easily seen that the *Arden-operation* mimics Arden's Lemma on the level of equations. To ensure the non-emptiness condition of Arden's Lemma we say that a right-hand side is *ardenable* provided

$$\text{ardenable } rhs \stackrel{\text{def}}{=} \forall Y r. (Y, r) \in rhs \longrightarrow \square \notin \mathcal{L}(r)$$

This allows us to prove a version of Arden's Lemma on the level of equations.

LEMMA 3.2. *Given an equation  $X = rhs$ . If  $X = \bigcup \mathcal{L} \text{ ' } rhs$ , *ardenable*  $rhs$ , and finite  $rhs$ , then  $X = \bigcup \mathcal{L} \text{ ' } (\text{Arden } X \text{ } rhs)$ .*

Our *ardenable* condition is slightly stronger than needed for applying Arden's Lemma, but we can still ensure that it holds throughout our algorithm of transforming equations into solved form.

The *substitution-operation* takes an equation of the form  $X = xrhs$  and substitutes it into the right-hand side  $rhs$ .

$$\begin{aligned} \text{Subst } rhs \ X \ xrhs &\stackrel{\text{def}}{=} \text{let} \\ &\quad rhs' = rhs - \{(X, r) \mid (X, r) \in rhs\} \\ &\quad r' = \dagger \{r \mid (X, r) \in rhs\} \\ &\text{in } rhs' \cup (xrhs \triangleleft r') \end{aligned}$$

We again delete first all occurrences of  $(X, r)$  in  $rhs$ ; we then calculate the regular expression corresponding to the deleted terms; finally we append this regular expression to  $xrhs$  and union it up with  $rhs'$ . When we use the substitution operation we will arrange it so that  $xrhs$  does not contain any occurrence of  $X$ . For example substituting the first equation in (6) into the right-hand side of the second, thus eliminating the equivalence class  $X_1$ , gives us the equation

$$X_2 = \lambda(TIMES \ ONE \ (ATOM \ c)) \quad (11)$$

With these two operations in place, we can define the operation that removes one equation from an equational systems  $ES$ . The operation *Subst\_all* substitutes an equation  $X = xrhs$  throughout an equational system  $ES$ ; *Remove* then completely removes such an equation from  $ES$  by substituting it to the rest of the equational system, but first eliminating all recursive occurrences of  $X$  by applying *Arden* to  $xrhs$ .

$$\begin{aligned} \text{Subst\_all } ES \ X \ xrhs &\stackrel{\text{def}}{=} \{(Y, \text{Subst } yrhs \ X \ xrhs) \mid (Y, yrhs) \in ES\} \\ \text{Remove } ES \ X \ xrhs &\stackrel{\text{def}}{=} \text{Subst\_all } (ES - \{(X, xrhs)\}) \ X \ (\text{Arden } X \ xrhs) \end{aligned}$$

Finally, we can define how an equational system should be solved. For this we will need to iterate the process of eliminating equations until only one equation will be left in the system. However, we do not just want to have any equation as being the last one, but the one involving the equivalence class for which we want to calculate the regular expression. Let us suppose this equivalence class is  $X$ . Since  $X$  is the one to be solved, in every iteration step we have to pick an equation to be eliminated that is different from  $X$ . In this way  $X$  is kept to the final step. The choice is implemented using Hilbert's choice operator, written *SOME* in the definition below.

$$\begin{aligned} \text{Iter } X \ ES &\stackrel{\text{def}}{=} \text{let} \\ &\quad (Y, yrhs) = \text{SOME } (Y, yrhs). (Y, yrhs) \in ES \wedge X \neq Y \\ &\text{in } \text{Remove } ES \ Y \ yrhs \end{aligned}$$

The last definition we need applies *Iter* over and over until a condition *Cond* is *not* satisfied anymore. This condition states that there are more than one equation left in the equational system  $ES$ . To solve an equational system we use Isabelle/HOL's *while*-operator as follows:

$$\text{Solve } X \ ES \stackrel{\text{def}}{=} \text{while } \text{Cond } (\text{Iter } X) \ ES$$

We are not concerned here with the definition of this operator (see [Berghofer and Nipkow 2002] for example), but note that we eliminate in each *Iter*-step a single equation, and therefore have a well-founded termination order by taking the cardinality of the equational system  $ES$ . This enables us to prove properties about our definition of *Solve* when we 'call' it with the equivalence class  $X$  and the initial equational system *Init* ( $UNIV // \approx_A$ ) from (9)

using the principle:

$$\begin{array}{l}
\text{invariant } (\text{Init } (\text{UNIV} // \approx_A)) \\
\forall ES. \text{invariant } ES \wedge \text{Cond } ES \longrightarrow \text{invariant } (\text{Iter } X \text{ } ES) \\
\forall ES. \text{invariant } ES \wedge \text{Cond } ES \longrightarrow \text{card } (\text{Iter } X \text{ } ES) < \text{card } ES \\
\forall ES. \text{invariant } ES \wedge \neg \text{Cond } ES \longrightarrow P \text{ } ES \\
\hline
P (\text{Solve } X (\text{Init } (\text{UNIV} // \approx_A)))
\end{array} \tag{12}$$

This principle states that given an invariant (which we will specify below) we can prove a property  $P$  involving  $\text{Solve}$ . For this we have to discharge the following proof obligations: first the initial equational system satisfies the invariant; second the iteration step  $\text{Iter}$  preserves the invariant as long as the condition  $\text{Cond}$  holds; third  $\text{Iter}$  decreases the termination order, and fourth that once the condition does not hold anymore then the property  $P$  must hold.

The property  $P$  in our proof will state that  $\text{Solve } X (\text{Init } (\text{UNIV} // \approx_A))$  returns with a single equation  $X = xrhs$  for some  $xrhs$ , and that this equational system still satisfies the invariant. In order to get the proof through, the invariant is composed of the following six properties:

$$\begin{array}{ll}
\text{invariant } ES & \stackrel{\text{def}}{=} \text{finite } ES & (\text{finiteness}) \\
& \wedge \forall (X, rhs) \in ES. \text{finite } rhs & (\text{finiteness } rhs) \\
& \wedge \forall (X, rhs) \in ES. X = \bigcup \mathcal{L} \text{ } rhs & (\text{soundness}) \\
& \wedge \forall X \text{ } rhs \text{ } rhs'. (X, rhs) \in ES \wedge (X, rhs') \in ES \longrightarrow rhs = rhs' & (\text{distinctness}) \\
& \wedge \forall (X, rhs) \in ES. \text{ardenable } rhs & (\text{ardenable}) \\
& \wedge \forall (X, rhs) \in ES. rhss \text{ } rhs \subseteq lhss \text{ } ES & (\text{validity})
\end{array}$$

The first two ensure that the equational system is always finite (number of equations and number of terms in each equation); the third makes sure the ‘meaning’ of the equations is preserved under our transformations. The other properties are a bit more technical, but are needed to get our proof through. Distinctness states that every equation in the system is distinct. *Ardenable* ensures that we can always apply the *Arden* operation. The last property states that every  $rhs$  can only contain equivalence classes for which there is an equation. Therefore  $lhss$  is just the set containing the first components of an equational system, while  $rhss$  collects all equivalence classes  $X$  in the terms of the form  $(X, r)$ . That means formally  $lhss \text{ } ES \stackrel{\text{def}}{=} \{X \mid (X, rhs) \in ES\}$  and  $rhss \text{ } rhs \stackrel{\text{def}}{=} \{X \mid (X, r) \in rhs\}$ .

It is straightforward to prove that the initial equational system satisfies the invariant.

LEMMA 3.3. *If finite (UNIV // ≈<sub>A</sub>) then invariant (Init (UNIV // ≈<sub>A</sub>)).*

PROOF. Finiteness is given by the assumption and the way how we set up the initial equational system. Soundness is proved in Lemma 3.1. Distinctness follows from the fact that the equivalence classes are disjoint. The *ardenable* property also follows from the setup of the initial equational system, as does validity.  $\square$

Next we show that  $\text{Iter}$  preserves the invariant.

LEMMA 3.4. *If invariant ES, (X, rhs) ∈ ES and Cond ES, then invariant (Iter X ES).*

PROOF. The argument boils down to choosing an equation  $Y = yrhs$  to be eliminated and to show that  $\text{Subst\_all } (ES - \{(Y, yrhs)\}) \text{ } Y (\text{Arden } Y \text{ } yrhs)$  preserves the invariant. We prove this as follows:

$\forall ES. \text{invariant } (ES \cup \{(Y, yrhs)\}) \text{ implies invariant } (Subst\_all \text{ } ES \text{ } Y \text{ } (Arden \text{ } Y \text{ } yrhs))$

Finiteness is straightforward, as the *Subst* and *Arden* operations keep the equational system finite. These operations also preserve soundness and distinctness (we proved soundness for *Arden* in Lemma 3.2). The property *ardenable* is clearly preserved because the append-operation cannot make a regular expression to match the empty string. Validity is given because *Arden* removes an equivalence class from *yrhs* and then *Subst\_all* removes *Y* from the equational system. Having proved the implication above, we can instantiate *ES* with  $ES - \{(Y, yrhs)\}$  which matches with our proof-obligation of *Subst\_all*. Since  $ES = ES - \{(Y, yrhs)\} \cup \{(Y, yrhs)\}$ , we can use the assumption to complete the proof.  $\square$

We also need the fact that *Iter* decreases the termination measure.

LEMMA 3.5. *If invariant ES,  $(X, rhs) \in ES$  and Cond ES, then  $card (Iter \text{ } X \text{ } ES) < card \text{ } ES$ .*

PROOF. By assumption we know that *ES* is finite and has more than one element. Therefore there must be an element  $(Y, yrhs) \in ES$  with  $(Y, yrhs) \neq (X, rhs)$ . Using the distinctness property we can infer that  $Y \neq X$ . We further know that *Remove ES Y yrhs* removes the equation  $Y = yrhs$  from the system, and therefore the cardinality of *Iter* strictly decreases.  $\square$

This brings us to our property we want to establish for *Solve*.

LEMMA 3.6. *If finite  $(UNIV // \approx_A)$  and  $X \in UNIV // \approx_A$  then there exists a rhs such that  $Solve \text{ } X \text{ } (Init (UNIV // \approx_A)) = \{(X, rhs)\}$  and invariant  $\{(X, rhs)\}$ .*

PROOF. In order to prove this lemma using (12), we have to use a slightly stronger invariant since Lemma 3.4 and 3.5 have the precondition that  $(X, rhs) \in ES$  for some *rhs*. This precondition is needed in order to choose in the *Iter*-step an equation that is not  $X = rhs$ . Therefore our invariant cannot be just *invariant ES*, but must be *invariant ES*  $\wedge$   $(\exists rhs. (X, rhs) \in ES)$ . By assumption  $X \in UNIV // \approx_A$  and Lemma 3.3, the more general invariant holds for the initial equational system. This is premise 1 of (12). Premise 2 is given by Lemma 3.4 and the fact that *Iter* might modify the *rhs* in the equation  $X = rhs$ , but does not remove it. Premise 3 of (12) is by Lemma 3.5. Now in premise 4 we like to show that there exists a *rhs* such that  $ES = \{(X, rhs)\}$  and that *invariant*  $\{(X, rhs)\}$  holds, provided the condition *Cond* does not hold. By the stronger invariant we know there exists such a *rhs* with  $(X, rhs) \in ES$ . Because *Cond* is not true, we know the cardinality of *ES* is 1. This means *ES* must actually be the set  $\{(X, rhs)\}$ , for which the invariant holds. This allows us to conclude that  $Solve \text{ } X \text{ } (Init (UNIV // \approx_A)) = \{(X, rhs)\}$  and *invariant*  $\{(X, rhs)\}$  hold, as needed.  $\square$

With this lemma in place we can show that for every equivalence class in  $UNIV // \approx_A$  there exists a regular expression.

LEMMA 3.7. *If finite  $(UNIV // \approx_A)$  and  $X \in UNIV // \approx_A$  then regular X.*

PROOF. By the preceding lemma, we know that there exists a *rhs* such that  $Solve \text{ } X \text{ } (Init (UNIV // \approx_A))$  returns the equation  $X = rhs$ , and that the invariant holds for this equation. That means we know  $X = \bigcup \mathcal{L} \text{ } rhs$ . We further know that this is equal to

$\bigcup \mathcal{L}^*(Arden X rhs)$  using the properties of the invariant and Lemma 3.2. Using the validity property for the equation  $X = rhs$ , we can infer that  $rhss rhs \subseteq \{X\}$  and because the *Arden* operation removes that  $X$  from  $rhs$ , that  $rhss (Arden X rhs) = \{\}$ . This means the right-hand side *Arden X rhs* can only consist of terms of the form  $\lambda(r)$ . So we can collect those (finitely many) regular expressions  $rs$  and have  $X = \mathcal{L}(\dagger rs)$ . With this we can conclude the proof.  $\square$

Lemma 3.7 allows us to finally give a proof for the first direction of the Myhill-Nerode Theorem.

PROOF OF THEOREM 3.1. By Lemma 3.7 we know that there exists a regular expression for every equivalence class in  $UNIV // \approx_A$ . Since *finals A* is a subset of  $UNIV // \approx_A$ , we also know that for every equivalence class in *finals A* there exists a regular expression. Moreover by assumption we know that *finals A* must be finite, and therefore there must be a finite set of regular expressions  $rs$  such that  $\bigcup \text{finals } A = \mathcal{L}(\dagger rs)$ . Since the left-hand side is equal to  $A$ , we can use  $\dagger rs$  as the regular expression that is needed in the theorem.  $\square$

Note that our algorithm for solving equational systems provides also a method for calculating a regular expression for the complement of a regular language: if we combine all regular expressions corresponding to equivalence classes not in *finals A*, then we obtain a regular expression for the complement language  $\bar{A}$ . This is similar to the usual construction of a ‘complement automaton’.

#### 4. MYHILL-NERODE, SECOND PART

In this section we will give a proof for establishing the second part of the Myhill-Nerode Theorem. It can be formulated in our setting as follows:

THEOREM 4.1. *Given  $r$  is a regular expression, then finite  $(UNIV // \approx_{\mathcal{L}(r)})$ .*

The proof will be by induction on the structure of  $r$ . It turns out the base cases are straightforward.

PROOF (BASE CASES). The cases for *ZERO*, *ONE* and *ATOM* are routine, because we can easily establish that

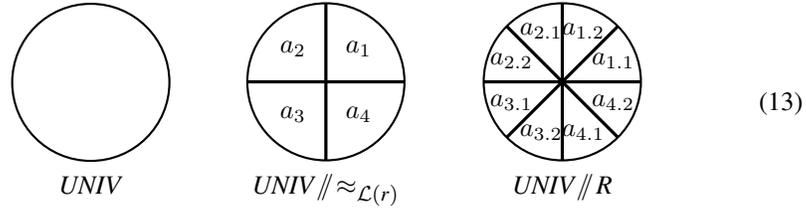
$$\begin{aligned} UNIV // \approx_{\{\}} &= \{UNIV\} \\ UNIV // \approx_{\{\square\}} &\subseteq \{\{\square\}, UNIV - \{\square\}\} \\ UNIV // \approx_{\{[c]\}} &\subseteq \{\{\square\}, \{[c]\}, UNIV - \{\square, [c]\}\} \end{aligned}$$

hold, which shows that  $UNIV // \approx_{\mathcal{L}(r)}$  must be finite.  $\square$

Much more interesting, however, are the inductive cases. They seem hard to be solved directly. The reader is invited to try.<sup>5</sup>

<sup>5</sup>The induction hypothesis is not strong enough to make any progress with the *TIME* and *STAR* cases.

In order to see how our proof proceeds consider the following suggestive picture given by Constable et al. [2000]:



The relation  $\approx_{\mathcal{L}(r)}$  partitions the set of all strings,  $UNIV$ , into some equivalence classes. To show that there are only finitely many of them, it suffices to show in each induction step that another relation, say  $R$ , has finitely many equivalence classes and refines  $\approx_{\mathcal{L}(r)}$ .

DEFINITION 4.1. A relation  $R_1$  refines  $R_2$  provided  $R_1 \subseteq R_2$ .

For constructing  $R$ , we will rely on some *tagging-functions* defined over strings. Given the inductive hypothesis, it will be easy to prove that the *range* of these tagging-functions is finite. The range of a function  $f$  is defined as

$$\text{range } f \stackrel{\text{def}}{=} f ' UNIV$$

that means we take the image of  $f$  w.r.t. all elements in the domain. With this we will be able to infer that the tagging-functions, seen as relations, give rise to finitely many equivalence classes. Finally we will show that the tagging-relations are more refined than  $\approx_{\mathcal{L}(r)}$ , which implies that  $UNIV // \approx_{\mathcal{L}(r)}$  must also be finite. We formally define the notion of a *tagging-relation* as follows.

DEFINITION 4.2 (TAGGING-RELATION). Given a tagging-function  $\text{tag}$ , then two strings  $x$  and  $y$  are tag-related provided

$$x \approx_{\text{tag}} y \stackrel{\text{def}}{=} \text{tag } x = \text{tag } y.$$

In order to establish finiteness of a set  $A$ , we shall use the following powerful principle from Isabelle/HOL's library.

$$\text{If finite } (f ' A) \text{ and inj\_on } f A \text{ then finite } A. \quad (14)$$

It states that if an image of a set under an injective function  $f$  (injective over this set) is finite, then the set  $A$  itself must be finite. We can use it to establish the following two lemmas.

LEMMA 4.1. If finite (range tag) then finite ( $UNIV // \approx_{\text{tag}}$ ).

PROOF. We set in (14),  $f$  to be  $X \mapsto \text{tag } ' X$ . We have  $\text{range } f$  to be a subset of  $\text{Pow}(\text{range } \text{tag})$ , which we know must be finite by assumption. Now  $f ' UNIV // \approx_{\text{tag}}$  is a subset of  $\text{range } f$ , and so also finite. Injectivity amounts to showing that  $X = Y$  under the assumptions that  $X, Y \in UNIV // \approx_{\text{tag}}$  and  $f X = f Y$ . From the assumptions we obtain  $x \in X$  and  $y \in Y$  with  $\text{tag } x = \text{tag } y$ . Since  $x$  and  $y$  are tag-related, this in turn means that the equivalence classes  $X$  and  $Y$  must be equal. Therefore (14) allows us to conclude with finite ( $UNIV // \approx_{\text{tag}}$ ).  $\square$

LEMMA 4.2. Given two equivalence relations  $R_1$  and  $R_2$ , whereby  $R_1$  refines  $R_2$ . If finite ( $UNIV // R_1$ ) then finite ( $UNIV // R_2$ ).

PROOF. We prove this lemma again using (14). This time we set  $f$  to be  $X \mapsto \{\llbracket x \rrbracket_{R_1} \mid x \in X\}$ . It is easy to see that  $\text{finite}(f \circ \text{UNIV} // R_2)$  because it is a subset of  $\text{Pow}(\text{UNIV} // R_1)$ , which must be finite by assumption. What remains to be shown is that  $f$  is injective on  $\text{UNIV} // R_2$ . This is equivalent to showing that two equivalence classes, say  $X$  and  $Y$ , in  $\text{UNIV} // R_2$  are equal, provided  $fX = fY$ . For  $X = Y$  to be equal, we have to find two elements  $x \in X$  and  $y \in Y$  such that they are  $R_2$  related. We know there exists a  $x \in X$  with  $X = \llbracket x \rrbracket_{R_2}$ . From the latter fact we can infer that  $\llbracket x \rrbracket_{R_1} \in fX$  and further  $\llbracket x \rrbracket_{R_1} \in fY$ . This means we can obtain a  $y$  such that  $\llbracket x \rrbracket_{R_1} = \llbracket y \rrbracket_{R_1}$  holds. Consequently  $x$  and  $y$  are  $R_1$ -related. Since by assumption  $R_1$  refines  $R_2$ , they must also be  $R_2$ -related, as we need to show.  $\square$

Chaining Lemma 4.1 and 4.2 together, means in order to show that  $\text{UNIV} // \approx_{\mathcal{L}(r)}$  is finite, we have to construct a tagging-function whose range can be shown to be finite and whose tagging-relation refines  $\approx_{\mathcal{L}(r)}$ . Let us attempt the *PLUS*-case first. We take as tagging-function

$$+tag A B x \stackrel{\text{def}}{=} (\llbracket x \rrbracket_{\approx_A}, \llbracket x \rrbracket_{\approx_B})$$

where  $A$  and  $B$  are some arbitrary languages. The reason for this choice is that we need to establish that  $\approx_{+tag A B}$  refines  $\approx_{A \cup B}$ . This amounts to showing  $x \approx_A y$  or  $x \approx_B y$  under the assumption  $x \approx_{+tag A B} y$ . As we shall see, this definition will provide us with just the right assumptions in order to get the proof through.

PROOF (*PLUS-CASE*). We can show in general, if  $\text{finite}(\text{UNIV} // \approx_A)$  and  $\text{finite}(\text{UNIV} // \approx_B)$  then  $\text{finite}(\text{UNIV} // \approx_A \times \text{UNIV} // \approx_B)$  holds. The range of  $+tag A B$  is a subset of this product set—so finite. For the refinement proof-obligation, we know that  $(\llbracket x \rrbracket_{\approx_A}, \llbracket x \rrbracket_{\approx_B}) = (\llbracket y \rrbracket_{\approx_A}, \llbracket y \rrbracket_{\approx_B})$  holds by assumption. Then clearly either  $x \approx_A y$  or  $x \approx_B y$ , as we needed to show. Finally we can discharge this case by setting  $A$  to  $\mathcal{L}(r_1)$  and  $B$  to  $\mathcal{L}(r_2)$ .  $\square$

The *TIMES*-case is slightly more complicated. We first prove the following lemma, which will aid the proof about refinement.

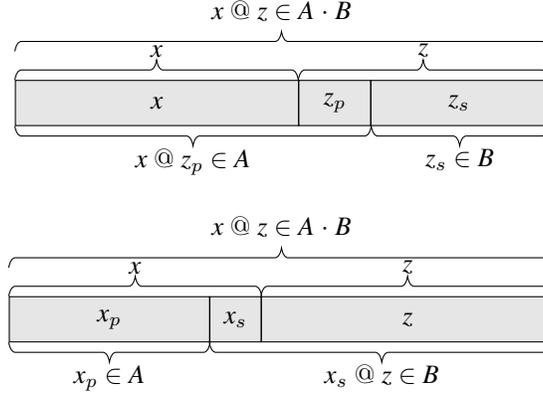
LEMMA 4.3. *The relation  $\approx_{tag}$  refines  $\approx_A$ , provided for all strings  $x, y$  and  $z$  we have that  $x \approx_{tag} y$  and  $x @ z \in A$  imply  $y @ z \in A$ .*

We therefore can analyse how the strings  $x @ z$  are in the language  $A$  and then construct an appropriate tagging-function to infer that  $y @ z$  are also in  $A$ . For this we will use the notion of the set of all possible *partitions* of a string:

$$\text{Partitions } x \stackrel{\text{def}}{=} \{(x_p, x_s) \mid x_p @ x_s = x\} \quad (15)$$

If we know that  $(x_p, x_s) \in \text{Partitions } x$ , we will refer to  $x_p$  as the *prefix* of the string  $x$ , and respectively to  $x_s$  as the *suffix*.

Now assuming  $x @ z \in A \cdot B$ , there are only two possible ways of how to ‘split’ this string to be in  $A \cdot B$ :



Either  $x$  and a prefix of  $z$  is in  $A$  and the rest in  $B$  (first picture) or there is a prefix of  $x$  in  $A$  and the rest is in  $B$  (second picture). In both cases we have to show that  $y @ z \in A \cdot B$ . The first case we will only go through if we know that  $x \approx_A y$  holds (\*). Because then we can infer from  $x @ z_p \in A$  that  $y @ z_p \in A$  holds for all  $z_p$ . In the second case we only know that  $x_p$  and  $x_s$  is one possible partition of the string  $x$ . We have to know that both  $x_p$  and the corresponding partition  $y_p$  are in  $A$ , and that  $x_s$  is ‘ $B$ -related’ to  $y_s$  (\*\*). From the latter fact we can infer that  $y_s @ z \in B$ . This will solve the second case. Taking the two requirements, (\*) and (\*\*), together we define the tagging-function in the *TIMES*-case as:

$$\times tag A B \stackrel{def}{=} ([x]_{\approx_A}, \{[x_s]_{\approx_B} \mid x_p \in A \wedge (x_p, x_s) \in Partitions x\})$$

Note that we have to make the assumption for all suffixes  $x_s$ , since we do not know anything about how the string  $x$  is partitioned. With this definition in place, let us prove the *TIMES*-case.

**PROOF (*TIMES*-CASE).** If *finite* ( $UNIV // \approx_A$ ) and *finite* ( $UNIV // \approx_B$ ) then *finite* ( $UNIV // \approx_A \times Pow (UNIV // \approx_B)$ ) holds. The range of  $\times tag A B$  is a subset of this product set, and therefore finite. For the refinement of  $\approx_A \cdot B$  and  $\approx_{\times tag A B}$ , we have by Lemma 4.3

$$\times tag A B x = \times tag A B y$$

and  $x @ z \in A \cdot B$ , and have to establish  $y @ z \in A \cdot B$ . As shown in the pictures above, there are two cases to be considered. First, there exists a  $z_p$  and  $z_s$  such that  $x @ z_p \in A$  and  $z_s \in B$ . By the assumption about  $\times tag A B$  we have  $[x]_{\approx_A} = [y]_{\approx_A}$  and thus  $x \approx_A y$ . Hence by the Myhill-Nerode Relation  $y @ z_p \in A$  holds. Using  $z_s \in B$ , we can conclude in this case with  $y @ z \in A \cdot B$  (recall  $z_p @ z_s = z$ ).

Second there exists a partition  $x_p$  and  $x_s$  with  $x_p \in A$  and  $x_s @ z \in B$ . We therefore have

$$[x_s]_{\approx_B} \in \{[x_s]_{\approx_B} \mid x_p \in A \wedge (x_p, x_s) \in Partitions x\}$$

and by the assumption about  $\times tag A B$  also

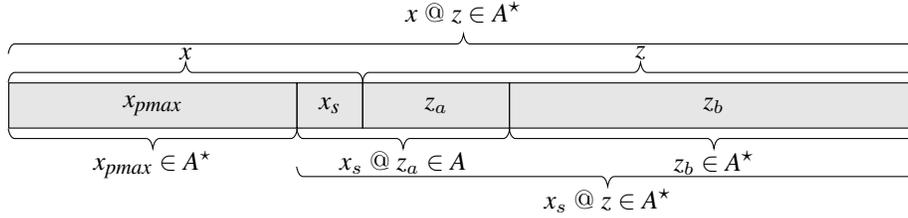
$$[x_s]_{\approx_B} \in \{[y_s]_{\approx_B} \mid y_p \in A \wedge (y_p, y_s) \in Partitions y\}$$

This means there must be a partition  $y_p$  and  $y_s$  such that  $y_p \in A$  and  $[x_s]_{\approx_B} = [y_s]_{\approx_B}$ . Unfolding the Myhill-Nerode Relation and together with the facts that  $x_p \in A$  and  $x_s @ z \in B$ , we obtain  $y_p \in A$  and  $y_s @ z \in B$ , as needed in this case. We again can complete the *TIMES*-case by setting  $A$  to  $\mathcal{L}(r_1)$  and  $B$  to  $\mathcal{L}(r_2)$ .  $\square$

The case for *STAR* is similar to *TIMES*, but poses a few extra challenges. To deal with them, we define first the notion of a *string prefix* and a *strict string prefix*:

$$\begin{aligned} x \leq y &\stackrel{\text{def}}{=} \exists z. y = x @ z \\ x < y &\stackrel{\text{def}}{=} x \leq y \wedge x \neq y \end{aligned}$$

When analysing the case of  $x @ z \in A^*$  and  $x$  is not the empty string, we have the following picture:



We can find a strict prefix  $x_p$  of  $x$  such that  $x_p \in A^*$ ,  $x_p < x$  and the rest  $x_s @ z \in A^*$ . For example the empty string  $\square$  would do (recall  $x \neq \square$ ). There are potentially many such prefixes, but there can only be finitely many of them (the string  $x$  is finite). Let us therefore choose the longest one and call it  $x_{pmax}$ . Now for the rest of the string  $x_s @ z$  we know it is in  $A^*$  and cannot be the empty string. By Property 2.1(iv), we can separate this string into two parts, say  $a$  and  $b$ , such that  $a \neq \square$ ,  $a \in A$  and  $b \in A^*$ . Now  $a$  must be strictly longer than  $x_s$ , otherwise  $x_{pmax}$  is not the longest prefix. That means  $a$  ‘overlaps’ with  $z$ , splitting it into two components  $z_a$  and  $z_b$ . For this we know that  $x_s @ z_a \in A$  and  $z_b \in A^*$ . To cut a story short, we have divided  $x @ z \in A^*$  such that we have a string  $a$  with  $a \in A$  that lies just on the ‘border’ of  $x$  and  $z$ . This string is  $x_s @ z_a$ .

In order to show that  $x @ z \in A^*$  implies  $y @ z \in A^*$ , we use the following tagging-function:

$$\star tag A x \stackrel{\text{def}}{=} \{ \llbracket x_s \rrbracket_{\approx A} \mid x_p < x \wedge x_p \in A^* \wedge (x_p, x_s) \in Partitions x \}$$

PROOF (*STAR*-CASE). If *finite* ( $UNIV // \approx_A$ ) then *finite* ( $Pow (UNIV // \approx_A)$ ) holds. The range of  $\star tag A$  is a subset of this set, and therefore finite. Again we have to show under the assumption  $x \approx_{\star tag A} y$  that  $x @ z \in A^*$  implies  $y @ z \in A^*$ .

We first need to consider the case that  $x$  is the empty string. From the assumption about strict prefixes in  $\approx_{\star tag A}$ , we can infer  $y$  is the empty string and then clearly have  $y @ z \in A^*$ . In case  $x$  is not the empty string, we can divide the string  $x @ z$  as shown in the picture above. By the tagging-function and the facts  $x_{pmax} \in A^*$  and  $x_{pmax} < x$ , we have

$$\llbracket x_s \rrbracket_{\approx A} \in \{ \llbracket x_s \rrbracket_{\approx A} \mid x_{pmax} < x \wedge x_{pmax} \in A^* \wedge (x_{pmax}, x_s) \in Partitions x \}$$

which by assumption is equal to

$$\llbracket x_s \rrbracket_{\approx A} \in \{ \llbracket y_s \rrbracket_{\approx A} \mid y_p < y \wedge y_p \in A^* \wedge (y_p, y_s) \in Partitions y \}$$

From this we know there exist a partition  $y_p$  and  $y_s$  with  $y_p \in A^*$  and also  $x_s \approx_A y_s$ . Unfolding the Myhill-Nerode Relation we know  $y_s @ z_a \in A$ . We also know that  $z_b \in A^*$ . Therefore  $y_p @ (y_s @ z_a) @ z_b \in A^*$ , which means  $y @ z \in A^*$ . The last step is to set  $A$  to  $\mathcal{L}(r)$  and thus complete the proof.  $\square$

## 5. SECOND PART PROVED USING PARTIAL DERIVATIVES

As we have seen in the previous section, in order to establish the second direction of the Myhill-Nerode Theorem, it is sufficient to find a more refined relation than  $\approx_{\mathcal{L}(r)}$  for which we can show that there are only finitely many equivalence classes. So far we showed this directly by induction on  $r$  using tagging-functions. However, there is also an indirect method to come up with such a refined relation by using derivatives of regular expressions introduced by Brzozowski [1964].

Assume the following two definitions for the *left-quotient* of a language, which we write as  $Der\ c\ A$  and  $Ders\ s\ A$  where  $c$  is a character and  $s$  a string, respectively:

$$\begin{aligned} Der\ c\ A &\stackrel{def}{=} \{s \mid [c] @ s \in A\} \\ Ders\ s\ A &\stackrel{def}{=} \{s' \mid s @ s' \in A\} \end{aligned}$$

In order to aid readability, we shall make use of the following abbreviation

$$Der\ s\ s\ A \stackrel{def}{=} \bigcup Ders\ s'\ A$$

where we apply the left-quotient to a set of languages and then combine the results. Clearly we have the following equivalence between the Myhill-Nerode Relation (Definition 3.1) and left-quotients

$$x \approx_A y \quad \text{if and only if} \quad Der\ x\ A = Der\ y\ A \quad (16)$$

It is also straightforward to establish the following properties of left-quotients

$$\begin{aligned} Der\ a\ \{\} &= \{\} \\ Der\ a\ \{\square\} &= \{\} \\ Der\ a\ \{[b]\} &= \text{if } a = b \text{ then } \{\square\} \text{ else } \{\} \\ Der\ a\ (A \cup B) &= Der\ a\ A \cup Der\ a\ B \\ Der\ c\ (A \cdot B) &= (Der\ c\ A) \cdot B \cup (\text{if } \square \in A \text{ then } Der\ c\ B \text{ else } \{\}) \\ Der\ c\ (A^*) &= (Der\ c\ A) \cdot A^* \\ Ders\ \square\ A &= A \\ Ders\ (c :: s)\ A &= Ders\ s\ (Der\ c\ A) \end{aligned} \quad (17)$$

Note that in the last equation we use the list-cons operator written  $_ :: _$ . The only interesting case is the case of  $A^*$  where we use Property 2.1(i) in order to infer that  $Der\ c\ (A^*) = Der\ c\ (A \cdot A^*)$ . We can then complete the proof by using the fifth equation and noting that  $Der\ c\ (A^*) \subseteq (Der\ c\ A) \cdot A^*$  provided  $\square \in A$ .

Brzozowski [1964] observed that the left-quotients for languages of regular expressions can be calculated directly using the notion of *derivatives of a regular expression*. We define this notion in Isabelle/HOL as follows:

$$\begin{aligned} der\ c\ (ZERO) &\stackrel{def}{=} ZERO \\ der\ c\ (ONE) &\stackrel{def}{=} ZERO \\ der\ c\ (ATOM\ d) &\stackrel{def}{=} \text{if } c = d \text{ then } ONE \text{ else } ZERO \\ der\ c\ (PLUS\ r_1\ r_2) &\stackrel{def}{=} PLUS\ (der\ c\ r_1)\ (der\ c\ r_2) \\ der\ c\ (TIMES\ r_1\ r_2) &\stackrel{def}{=} \text{if } \delta(r_1) \text{ then } PLUS\ (TIMES\ (der\ c\ r_1)\ r_2)\ (der\ c\ r_2) \\ &\quad \text{else } TIMES\ (der\ c\ r_1)\ r_2 \end{aligned}$$

$$\begin{aligned}
\text{der } c \text{ (STAR } r) & \stackrel{\text{def}}{=} \text{TIMES (der } c \text{ } r) \text{ (STAR } r) \\
\text{ders } [] \text{ } r & \stackrel{\text{def}}{=} r \\
\text{ders } (c :: s) \text{ } r & \stackrel{\text{def}}{=} \text{ders } s \text{ (der } c \text{ } r)
\end{aligned}$$

The last two clauses extend derivatives from characters to strings. The boolean function  $\delta(r)$  needed in the *TIMES*-case tests whether a regular expression can recognise the empty string. It can be defined as follows.

$$\begin{aligned}
\delta(\text{ZERO}) & \stackrel{\text{def}}{=} \text{False} \\
\delta(\text{ONE}) & \stackrel{\text{def}}{=} \text{True} \\
\delta(\text{ATOM } c) & \stackrel{\text{def}}{=} \text{False} \\
\delta(\text{PLUS } r_1 \text{ } r_2) & \stackrel{\text{def}}{=} \delta(r_1) \vee \delta(r_2) \\
\delta(\text{TIMES } r_1 \text{ } r_2) & \stackrel{\text{def}}{=} \delta(r_1) \wedge \delta(r_2) \\
\delta(\text{STAR } r) & \stackrel{\text{def}}{=} \text{True}
\end{aligned}$$

By induction on the regular expression  $r$ , respectively on the string  $s$ , one can easily show that left-quotients and derivatives of regular expressions relate as follows (see for example [Sakarovitch 2009]):

$$\begin{aligned}
\text{Der } c \text{ (}\mathcal{L}(r)\text{)} &= \mathcal{L}(\text{der } c \text{ } r) \\
\text{Ders } s \text{ (}\mathcal{L}(r)\text{)} &= \mathcal{L}(\text{ders } s \text{ } r)
\end{aligned} \tag{18}$$

The importance of this fact in the context of the Myhill-Nerode Theorem is that we can use (16) and (18) in order to establish that

$$x \approx_{\mathcal{L}(r)} y \quad \text{if and only if} \quad \mathcal{L}(\text{ders } x \text{ } r) = \mathcal{L}(\text{ders } y \text{ } r).$$

holds and hence

$$x \approx_{\mathcal{L}(r)} y \quad \text{provided} \quad \text{ders } x \text{ } r = \text{ders } y \text{ } r \tag{19}$$

This means the right-hand side (seen as a relation) refines the Myhill-Nerode Relation. Consequently, we can use  $\approx_{(\lambda x. \text{ders } x \text{ } r)}$  as a tagging-relation. However, in order to be useful for the second part of the Myhill-Nerode Theorem, we have to be able to establish that for the corresponding language there are only finitely many derivatives—thus ensuring that there are only finitely many equivalence classes. Unfortunately, this is not true in general. Sakarovitch gives an example where a regular expression has infinitely many derivatives w.r.t. the language  $(ab)^* \cup (ab)^*a$ , which is formally written in our notation as  $\{[a,b]^*\} \cup (\{[a,b]^*\} \cdot \{[a]\})$  (see [Sakarovitch 2009, Page 141]).

What Brzozowski [1964] established is that for every language there *are* only finitely ‘dissimilar’ derivatives for a regular expression. Two regular expressions are said to be

*similar* provided they can be identified using the using the *ACI*-identities:

$$\begin{aligned}
(A) \quad & PLUS (PLUS r_1 r_2) r_3 \equiv PLUS r_1 (PLUS r_2 r_3) \\
(C) \quad & PLUS r_1 r_2 \equiv PLUS r_2 r_1 \\
(I) \quad & PLUS r r \equiv r
\end{aligned} \tag{20}$$

Carrying this idea through, we must not consider the set of all derivatives, but the one modulo *ACI*. In principle, this can be done formally, but it is very painful in a theorem prover (since there is no direct characterisation of the set of dissimilar derivatives).

Fortunately, there is a much simpler approach using *partial derivatives*. They were introduced by Antimirov [1995] and can be defined in Isabelle/HOL as follows:

$$\begin{aligned}
pder\ c\ (ZERO) & \stackrel{def}{=} \{\} \\
pder\ c\ (ONE) & \stackrel{def}{=} \{\} \\
pder\ c\ (ATOM\ d) & \stackrel{def}{=} \text{if } c = d \text{ then } \{ONE\} \text{ else } \{\} \\
pder\ c\ (PLUS\ r_1\ r_2) & \stackrel{def}{=} pder\ c\ r_1 \cup pder\ c\ r_2 \\
pder\ c\ (TIMES\ r_1\ r_2) & \stackrel{def}{=} \text{if } \delta(r_1) \text{ then } TIMESS\ (pder\ c\ r_1)\ r_2 \cup pder\ c\ r_2 \\
& \quad \text{else } TIMESS\ (pder\ c\ r_1)\ r_2 \\
pder\ c\ (STAR\ r) & \stackrel{def}{=} TIMESS\ (pder\ c\ r)\ (STAR\ r) \\
pders\ []\ r & \stackrel{def}{=} \{r\} \\
pders\ (c :: s)\ r & \stackrel{def}{=} \bigcup (pders\ s)\ ' (pder\ c\ r)
\end{aligned}$$

Again the last two clauses extend partial derivatives from characters to strings. Unlike ‘simple’ derivatives, the functions for partial derivatives return sets of regular expressions. In the *TIMES* and *STAR* cases we therefore use the auxiliary definition

$$TIMESS\ rs\ r \stackrel{def}{=} \{TIMES\ r'\ r \mid r' \in rs\}$$

in order to ‘sequence’ a regular expression with a set of regular expressions. Note that in the last clause we first build the set of partial derivatives w.r.t the character  $c$ , then build the image of this set under the function  $pders\ s$  and finally ‘union up’ all resulting sets. It will be convenient to introduce for this the following abbreviation

$$pderss\ s\ rs \stackrel{def}{=} \bigcup pders\ s\ ' rs$$

which simplifies the last clause of  $pders$  to

$$pders\ (c :: s)\ r \stackrel{def}{=} pderss\ s\ (pder\ c\ r)$$

Partial derivatives can be seen as having the *ACI*-identities already built in: taking the partial derivatives of the regular expressions in (20) gives us in each case equal sets. Antimirov [1995] showed a similar result to (18) for partial derivatives, namely

$$\begin{aligned}
(i) \quad & Der\ c\ (\mathcal{L}(r)) = \bigcup \mathcal{L}\ ' pder\ c\ r \\
(ii) \quad & Ders\ s\ (\mathcal{L}(r)) = \bigcup \mathcal{L}\ ' pders\ s\ r
\end{aligned} \tag{21}$$

PROOF. The first fact is by a simple induction on  $r$ . For the second we slightly modify Antimirov's proof by performing an induction on  $s$  where we generalise over all  $r$ . That means in the *cons*-case the induction hypothesis is

$$(IH) \quad \forall r. Ders\ s\ (\mathcal{L}(r)) = \bigcup \mathcal{L}'\ pders\ s\ r$$

With this we can establish

$$\begin{aligned} Ders\ (c :: s)\ (\mathcal{L}(r)) &= Ders\ s\ (Der\ c\ (\mathcal{L}(r))) && \text{by def.} \\ &= Ders\ s\ (\bigcup \mathcal{L}'\ pder\ c\ r) && \text{by (21.i)} \\ &= Derss\ s\ (\mathcal{L}'\ pder\ c\ r) && \text{by def. of } Ders \\ &= \bigcup \mathcal{L}'\ pderss\ s\ (pder\ c\ r) && \text{by IH} \\ &= \bigcup \mathcal{L}'\ pders\ (c :: s)\ r && \text{by def.} \end{aligned}$$

Note that in order to apply the induction hypothesis in the fourth equation, we need the generalisation over all regular expressions  $r$ . The case for the empty string is routine and omitted.  $\square$

Taking (18) and (21) together gives the relationship between languages of derivatives and partial derivatives

$$\begin{aligned} (i) \quad \mathcal{L}(der\ c\ r) &= \bigcup \mathcal{L}'\ pder\ c\ r \\ (ii) \quad \mathcal{L}(ders\ s\ r) &= \bigcup \mathcal{L}'\ pders\ s\ r \end{aligned} \quad (22)$$

These two properties confirm the observation made earlier that by using sets, partial derivatives have the *ACI*-identities of derivatives already built in.

Antimirov also proved that for every language and every regular expression there are only finitely many partial derivatives, whereby the set of partial derivatives of  $r$  w.r.t. a language  $A$  is defined as

$$pdersl\ A\ r \stackrel{\text{def}}{=} \bigcup_{x \in A} pders\ x\ r \quad (23)$$

THEOREM 5.1 [ANTIMIROV 1995]. *For every language  $A$  and every regular expression  $r$ , finite ( $pdersl\ A\ r$ ).*

Antimirov's proof first establishes this theorem for the language  $UNIV^+$ , which is the set of all non-empty strings. For this he proves:

$$\begin{aligned} pdersl\ UNIV^+\ (ZERO) &= \{\} \\ pdersl\ UNIV^+\ (ONE) &= \{\} \\ pdersl\ UNIV^+\ (ATOM\ c) &= \{ONE\} \\ pdersl\ UNIV^+\ (PLUS\ r_1\ r_2) &= pdersl\ UNIV^+\ r_1 \cup pdersl\ UNIV^+\ r_2 \\ pdersl\ UNIV^+\ (TIMES\ r_1\ r_2) &\subseteq TIMESS\ (pdersl\ UNIV^+\ r_1)\ r_2 \cup pdersl\ UNIV^+\ r_2 \\ pdersl\ UNIV^+\ (STAR\ r) &\subseteq TIMESS\ (pdersl\ UNIV^+\ r)\ (STAR\ r) \end{aligned} \quad (24)$$

from which one can deduce by induction on  $r$  that

$$finite\ (pdersl\ UNIV^+\ r)$$

holds. Now Antimirov's theorem follows because

$$pdersl\ UNIV\ r = pders\ \square\ r \cup pdersl\ UNIV^+\ r$$

and for all languages  $A$ ,  $pdersl\ A\ r$  is a subset of  $pdersl\ UNIV\ r$ . Since we follow Antimirov's proof quite closely in our formalisation (only the last two cases of (24) involve some non-routine induction arguments), we omit the details.

Let us now return to our proof for the second direction in the Myhill-Nerode Theorem. The point of the above calculations is to use  $\approx_{(\lambda x. pders\ x\ r)}$  as tagging-relation.

PROOF OF THEOREM 4.1 (SECOND VERSION). Using (16) and (21) we can easily infer that

$$x \approx_{\mathcal{L}(r)} y \quad \text{provided} \quad pders\ x\ r = pders\ y\ r$$

which means the tagging-relation  $\approx_{(\lambda x. pders\ x\ r)}$  refines  $\approx_{\mathcal{L}(r)}$ . So we know by Lemma 4.2,  $finite\ (UNIV // \approx_{\mathcal{L}(r)})$  holds if  $finite\ (UNIV // \approx_{(\lambda x. pders\ x\ r)})$ . In order to establish the latter, we can use Lemma 4.1 and show that the range of the tagging-function  $\lambda x. pders\ x\ r$  is finite. For this recall Definition 23, which gives us that

$$pdersl\ UNIV\ r \stackrel{def}{=} \bigcup_x pders\ x\ r$$

Now the range of  $\lambda x. pders\ x\ r$  is a subset of  $Pow\ (pdersl\ UNIV\ r)$ , which we know is finite by Theorem 5.1. Consequently there are only finitely many equivalence classes of  $\approx_{(\lambda x. pders\ x\ r)}$ . This relation refines  $\approx_{\mathcal{L}(r)}$ , and therefore we can again conclude the second part of the Myhill-Nerode Theorem.  $\square$

## 6. CLOSURE PROPERTIES OF REGULAR LANGUAGES

The beauty of regular languages is that they are closed under many set operations. Closure under union, concatenation and Kleene-star are trivial to establish given our definition of regularity (recall Definition 1.2). More interesting in our setting is the closure under complement, because it seems difficult to construct a regular expression for the complement language by direct means. However the existence of such a regular expression can now be easily proved using both parts of the Myhill-Nerode Theorem, since

$$s_1 \approx_A s_2 \text{ if and only if } s_1 \approx_{\bar{A}} s_2$$

holds for any strings  $s_1$  and  $s_2$ . Therefore  $A$  and the complement language  $\bar{A}$  give rise to the same partitions. So if one is finite, the other is too, and vice versa. As noted earlier, our algorithm for solving equational systems actually calculates a regular expression for the complement language. Calculating such a regular expression via automata using the standard method would be quite involved. It includes the steps: regular expression  $\Rightarrow$  non-deterministic automaton  $\Rightarrow$  deterministic automaton  $\Rightarrow$  complement automaton  $\Rightarrow$  regular expression. Clearly not something you want to formalise in a theorem prover if it is cumbersome to reason about automata.

A perhaps surprising fact is that regular languages are closed under any left-quotient. Define

$$Dersl\ B\ A \stackrel{def}{=} \bigcup_{x \in B} Ders\ x\ A$$

and assume  $B$  is any language and  $A$  is regular, then  $Dersl\ B\ A$  is regular. To see this consider the following argument using partial derivatives (which we used in Section 5): From  $A$  being regular we know there exists a regular expression  $r$  such that  $A = \mathcal{L}(r)$ . We also know that  $pdersl\ B\ r$  is finite for every language  $B$  and regular expression  $r$  (recall Theorem 5.1). By definition and (21) we have

$$Dersl\ B\ (\mathcal{L}(r)) = \bigcup \mathcal{L}\ ' pdersl\ B\ r \tag{25}$$

Since there are only finitely many regular expressions in  $pdersl B r$ , we know by (2) that there exists a regular expression so that the right-hand side of (25) is equal to the language  $\mathcal{L}(+(pdersl B r))$ . Thus the regular expression  $+(pdersl B r)$  verifies that  $Dersl B A$  is regular.

Even more surprising is the fact given first by Haines [1969] stating that for every language  $A$ , the language consisting of all (scattered) substrings of  $A$  is regular (see also [Shallit 2008; Fenner et al. 2009]). A (scattered) substring can be obtained by striking out zero or more characters from a string. This can be defined inductively in Isabelle/HOL by the following three rules:

$$\frac{}{\boxed{\phantom{x}} \preceq y} \quad \frac{x \preceq y}{x \preceq c :: y} \quad \frac{x \preceq y}{c :: x \preceq c :: y}$$

It is straightforward to prove that  $\preceq$  is a partial order. Now define the *language of substrings* and *superstrings* of a language  $A$  respectively as

$$\begin{aligned} Sub A &\stackrel{def}{=} \{x \mid \exists y \in A. x \preceq y\} \\ Sup A &\stackrel{def}{=} \{x \mid \exists y \in A. y \preceq x\} \end{aligned}$$

We like to establish

**THEOREM 6.1 [HAINES 1969].** *For every language  $A$ , the languages (i)  $Sub A$  and (ii)  $Sup A$  are regular.*

Our proof follows the one given by Shallit [2008, Pages 92–95], except that we use Higman’s Lemma, which is already proved in the Isabelle/HOL library by Sternagel. Higman’s Lemma allows us to infer that every language  $A$  of antichains, satisfying

$$\forall x, y \in A. x \neq y \longrightarrow x \not\preceq y \wedge y \not\preceq x \tag{26}$$

is finite.

The first step in our proof of Theorem 6.1 is to establish the following simple properties for  $Sup$

$$\begin{aligned} Sup \{\} &\stackrel{def}{=} \{\} \\ Sup \{\boxed{\phantom{x}}\} &\stackrel{def}{=} UNIV \\ Sup \{[c]\} &\stackrel{def}{=} UNIV \cdot \{[c]\} \cdot UNIV \\ Sup (A \cup B) &\stackrel{def}{=} Sup A \cup Sup B \\ Sup (A \cdot B) &\stackrel{def}{=} Sup A \cdot Sup B \\ Sup (A^*) &\stackrel{def}{=} UNIV \end{aligned} \tag{27}$$

whereby the last equation follows from the fact that  $A^*$  contains the empty string. With these properties at our disposal we can establish the lemma

**LEMMA 6.1.** *If  $A$  is regular, then also  $Sup A$ .*

**PROOF.** Since our alphabet is finite, we have a regular expression, written  $ALL$ , that matches every string. Using this regular expression we can inductively define the operation  $r \uparrow$

$$\begin{aligned}
(ZERO)\uparrow &\stackrel{def}{=} ZERO \\
(ONE)\uparrow &\stackrel{def}{=} ALL \\
(ATOM\ c)\uparrow &\stackrel{def}{=} TIMES\ ALL\ (TIMES\ (ATOM\ c)\ ALL) \\
(PLUS\ r_1\ r_2)\uparrow &\stackrel{def}{=} PLUS\ (r_1)\uparrow\ (r_2)\uparrow \\
(TIMES\ r_1\ r_2)\uparrow &\stackrel{def}{=} TIMES\ (r_1)\uparrow\ (r_2)\uparrow \\
(STAR\ r)\uparrow &\stackrel{def}{=} ALL
\end{aligned}$$

and use (27) to establish that  $\mathcal{L}((r)\uparrow) = Sup(\mathcal{L}(r))$  holds. This shows that  $Sup\ A$  is regular, provided  $A$  is.  $\square$

Now we can prove the main lemma w.r.t.  $Sup$ , namely

LEMMA 6.2. *For every language  $A$ , there exists a finite language  $M$  such that*

$$Sup\ M = Sup\ A .$$

PROOF. For  $M$  we take the set of all minimal elements of  $A$ . An element  $x$  is said to be *minimal* in  $A$  provided

$$min_A\ x \stackrel{def}{=} \forall y \in A. y \preceq x \longrightarrow x \preceq y$$

By Higman's Lemma (26) we know that  $M \stackrel{def}{=} \{x \in A \mid min_A\ x\}$  is finite, since every minimal element is incomparable, except with itself. It is also straightforward to show that  $Sup\ M \subseteq Sup\ A$ . For the other direction we have  $x \in Sup\ A$ . From this we obtain a  $y$  such that  $y \in A$  and  $y \preceq x$ . Since we have that the relation  $\{(y, x) \mid y \preceq x \wedge x \neq y\}$  is well-founded, there must be a minimal element  $z$  such that  $z \in A$  and  $z \preceq y$ , and hence by transitivity also  $z \preceq x$  (here we deviate from the argument given by Shallit [2008], because Isabelle/HOL provides already an extensive infrastructure for reasoning about well-foundedness). Since  $z$  is minimal and an element in  $A$ , we also know that  $z$  is in  $M$ . From this together with  $z \preceq x$ , we can infer that  $x$  is in  $Sup\ M$ , as required.  $\square$

This lemma allows us to establish the second part of Theorem 6.1.

PROOF OF THE SECOND PART OF THEOREM 6.1. Given any language  $A$ , by Lemma 6.2 we know there exists a finite, and thus regular, language  $M$ . We further have  $Sup\ M = Sup\ A$ , which establishes the second part.  $\square$

In order to establish the first part of this theorem, we use the property proved by Shallit [2008], namely that

$$\overline{Sub\ A} = Sup\ (\overline{Sub\ A}) \tag{28}$$

holds. Now the first part of Theorem 6.1 is a simple consequence of the second part.

PROOF OF THE FIRST PART OF THEOREM 6.1. By the second part, we know the right-hand side of (28) is regular, which means  $\overline{Sub\ A}$  is regular. But since we established already that regularity is preserved under complement (using Myhill-Nerode), also  $Sub\ A$  must be regular.  $\square$

Finally we like to show that the Myhill-Nerode Theorem is also convenient for establishing the non-regularity of languages. For this we use the following version of the Continuation Lemma (see for example [Rosenberg 2006]).

LEMMA 6.3 (CONTINUATION LEMMA). *If a language  $A$  is regular and a set of strings  $B$  is infinite, then there exist two distinct strings  $x$  and  $y$  in  $B$  such that  $x \approx_A y$ .*

This lemma can be easily deduced from the Myhill-Nerode Theorem and the Pigeonhole Principle: Since  $A$  is regular, there can be only finitely many equivalence classes. Hence an infinite set must contain at least two strings that are in the same equivalence class, that is they need to be related by the Myhill-Nerode Relation.

Using this lemma, it is straightforward to establish that the language  $A \stackrel{\text{def}}{=} \bigcup_n a^n @ b^n$  is not regular ( $a^n$  stands for the strings consisting of  $n$  times the character  $a$ ; similarly for  $b^n$ ). For this consider the infinite set  $B \stackrel{\text{def}}{=} \bigcup_n a^n$ .

LEMMA 6.4. *No two distinct strings in set  $B$  are Myhill-Nerode related by language  $A$ .*

PROOF. After unfolding the definition of  $B$ , we need to establish that given  $i \neq j$ , the strings  $a^i$  and  $a^j$  are not Myhill-Nerode related by  $A$ . That means we have to show that  $\forall z. a^i @ z \in A = a^j @ z \in A$  leads to a contradiction. Let us take  $b^i$  for  $z$ . Then we know  $a^i @ b^i \in A$ . But since  $i \neq j$ ,  $a^j @ b^i \notin A$ . Therefore  $a^i$  and  $a^j$  cannot be Myhill-Nerode related by  $A$ , and we are done.  $\square$

To conclude the proof of non-regularity for the language  $A$ , the Continuation Lemma and the lemma above lead to a contradiction assuming  $A$  is regular. Therefore the language  $A$  is not regular, as we wanted to show.

## 7. CONCLUSION AND RELATED WORK

In this paper we took the view that a regular language is one where there exists a regular expression that matches all of its strings. Regular expressions can conveniently be defined as a datatype in theorem provers. For us it was therefore interesting to find out how far we can push this point of view. We have established in Isabelle/HOL both directions of the Myhill-Nerode Theorem.

THEOREM 7.1 (MYHILL-NERODE THEOREM). *A language  $A$  is regular if and only if finite ( $UNIV // \approx_A$ ).*

Having formalised this theorem means we pushed our point of view quite far. Using this theorem we can obviously prove when a language is *not* regular—by establishing that it has infinitely many equivalence classes generated by the Myhill-Nerode Relation (this is usually the purpose of the Pumping Lemma). We can also use it to establish the standard textbook results about closure properties of regular languages. Interesting is the case of closure under complement, because it seems difficult to construct a regular expression for the complement language by direct means. However the existence of such a regular expression can be easily proved using the Myhill-Nerode Theorem.

While regular expressions are convenient, they have some limitations. One is that there seems to be no method of calculating a minimal regular expression (for example in terms of length) for a regular language, like there is for automata. On the other hand, efficient regular expression matching, without using automata, poses no problem as shown by Owens et al. [2009]. For an implementation of a simple regular expression matcher, whose correctness has been formally established, we refer the reader to Owens and Slind [2008]. In our opinion, their formalisation is considerably slicker than for example the approach to regular expression matching taken by Harper [1999] and by Yi [2006].

Our proof of the first direction is very much inspired by *Brzozowski's algebraic method* [1964] used to convert a finite automaton to a regular expression. The close connection can be seen by considering the equivalence classes as the states of the minimal automaton for the regular language. However there are some subtle differences. Because our equivalence classes (or correspondingly states) arise from the Myhill-Nerode Relation, the most natural choice is to characterise each state with the set of strings starting from the initial state leading up to that state. Usually, however, the states are characterised as the strings starting from that state leading to the terminal states. The first choice has consequences about how the initial equational system is set up. We have the  $\lambda$ -term on our 'initial state', while Brzozowski has it on the terminal states. This means we also need to reverse the direction of Arden's Lemma. We have not found anything in the 'pencil-and-paper-reasoning' literature about our way of proving the first direction of the Myhill-Nerode Theorem, but it appears to be folklore.

We presented two proofs for the second direction of the Myhill-Nerode Theorem. One direct proof using tagging-functions and another using partial derivatives. This part of our work is where our method using regular expressions shines, because we can completely side-step the standard argument (for example used by Kozen [1997]) where automata need to be composed. However, it is also the direction where we had to spend most of the 'conceptual' time, as our first proof based on tagging-functions is new for establishing the Myhill-Nerode Theorem. All standard proofs of this direction proceed by arguments over automata.

The indirect proof for the second direction arose from our interest in Brzozowski's derivatives for regular expression matching. While Brzozowski [1964] already established that there are only finitely many dissimilar derivatives for every regular expression, this result is not as straightforward to formalise in a theorem prover as one might wish. The reason is that the set of dissimilar derivatives is not defined inductively, but in terms of an ACI-equivalence relation. This difficulty prevented for example Krauss and Nipkow [2012] to prove termination of their equivalence checker for regular expressions. Their checker is based on Brzozowski's derivatives and for their argument the lack of a formal proof of termination is not crucial (it merely lets them "sleep better" [Krauss and Nipkow 2012]). We expect that their development simplifies by using partial derivatives, instead of derivatives, and that the termination of the algorithm can be formally established (the main ingredient is Theorem 5.1). However, since partial derivatives use sets of regular expressions, one needs to carefully analyse whether the resulting algorithm is still executable. Given the infrastructure for executable sets introduced by Haftmann [2009] in Isabelle/HOL, it should.

We started out by claiming that in a theorem prover it is easier to reason about regular expressions than about automata. Here are some numbers: Our formalisation of the Myhill-Nerode Theorem consists of 780 lines of Isabelle/Isar code for the first direction and 460 for the second (the one based on tagging-functions), plus around 300 lines of standard material about regular languages. The formalisation of derivatives and partial derivatives shown in Section 5 consists of 390 lines of code. The closure properties in Section 6 (except Theorem 6.1) can be established in 100 lines of code. The Continuation Lemma and the non-regularity of  $a^n b^n$  require 70 lines of code. The algorithm for solving equational systems, which we used in the first direction, is conceptually relatively simple. Still the use of sets over which the algorithm operates means it is not as easy to formalise as one might

wish. However, it seems sets cannot be avoided since the ‘input’ of the algorithm consists of equivalence classes and we cannot see how to reformulate the theory so that we can use lists or matrices. Lists would be much easier to reason about, since we can define functions over them by recursion. For sets we have to use set-comprehensions, which is slightly unwieldy. Matrices would allow us to use the slick formalisation by Nipkow [2011] of the Gauss-Jordan algorithm.

While our formalisation might appear large, it should be seen in the context of the work done by Constable et al. [2000] who formalised the Myhill-Nerode Theorem in Nuprl using automata. They write that their four-member team would need something on the magnitude of 18 months for their formalisation of the first eleven chapters of the textbook by Hopcroft and Ullman [1969], which includes the Myhill-Nerode theorem. It is hard to gauge the size of a formalisation in Nuprl, but from what is shown in the Nuprl Math Library about their development it seems *substantially* larger than ours. We attribute this to our use of regular expressions, which meant we did not need to ‘fight’ the theorem prover. Recently, Lammich and Tuerk [2012] formalised Hopcroft’s algorithm in Isabelle/HOL (in 7000 lines of code) using an automata library of 27000 lines of code. Also, Filliâtre [1997] reports that his formalisation in Coq of automata theory and Kleene’s theorem is “rather big”. Almeida et al. [2010] reported about another formalisation of regular languages in Coq. Their main result is the correctness of Mirkin’s construction of an automaton from a regular expression using partial derivatives. This took approximately 10600 lines of code. Braibant [2012] formalised a large part of regular language theory and Kleene algebras in Coq. While he is mainly interested in implementing decision procedures for Kleene algebras, his library includes a proof of the Myhill-Nerode theorem. He reckons that our “development is more concise” than his one based on matrices [Braibant 2012, Page 67]. He writes that there is no conceptual problems with formally reasoning about matrices for automata, but notes “intrinsic difficult[ies]” when working with matrices in Coq, which is the sort of ‘fighting’ one would encounter also in other theorem provers.

In terms of time, the estimate for our formalisation is that we needed approximately 3 months and this included the time to find our proof arguments. Unlike Constable et al. [2000], who were able to follow the Myhill-Nerode proof by Hopcroft and Ullman [1969], we had to find our own arguments. So for us the formalisation was not the bottleneck. The code of our formalisation [Wu et al. 2011b] can be found in the Archive of Formal Proofs at <http://afp.sourceforge.net/entries/Myhill-Nerode.shtml>.

**Acknowledgements:** We are grateful for the comments we received from Larry Paulson. Tobias Nipkow made us aware of the properties in Theorem 6.1 and Tjark Weber helped us with proving them. Christian Sternagel provided us with a version of Higman’s Lemma that applies to arbitrary, but finite alphabets.

## REFERENCES

- ALMEIDA, J. B., MORIERA, N., PEREIRA, D., AND DE SOUSA, S. M. 2010. Partial Derivative Automata Formalized in Coq. In *Proc. of the 15th International Conference on Implementation and Application of Automata*. LNCS, vol. 6482. 59–68.
- ANTIMIROV, V. 1995. Partial Derivatives of Regular Expressions and Finite Automata Constructions. *Theoretical Computer Science* 155, 291–319.
- ASPERTI, A. 2012. A Compact Proof of Decidability for Regular Expression Equivalence. In *Proc. of the 3rd International Conference on Interactive Theorem Proving*. LNCS, vol. 7406. 283–298.
- BERGHOFER, S. AND NIPKOW, T. 2002. Executing Higher Order Logic. In *Proc. of the International Workshop on Types for Proofs and Programs*. LNCS, vol. 2277. 24–40.

- BERGHOFER, S. AND REITER, M. 2009. Formalizing the Logic-Automaton Connection. In *Proc. of the 22nd International Conference on Theorem Proving in Higher Order Logics*. LNCS, vol. 5674. 147–163.
- BRAIBANT, T. 2012. Kleene Algebras, Rewriting Modulo AC, and Circuits in Coq. Ph.D. thesis, University of Grenoble.
- BRZOZOWSKI, J. A. 1964. Derivatives of Regular Expressions. *Journal of the ACM* 11, 4, 481–494.
- CHURCH, A. 1940. A Formulation of the Simple Theory of Types. *Journal of Symbolic Logic* 5, 2, 56–68.
- CONSTABLE, R. L., JACKSON, P. B., NAUMOV, P., AND URIBE, J. C. 2000. Constructively Formalizing Automata Theory. In *Proof, Language, and Interaction*. MIT Press, 213–238.
- COQUAND, T. AND SILES, V. 2011. A Decision Procedure for Regular Expression Equivalence in Type Theory. In *Proc. of the 1st Conference on Certified Programs and Proofs*. LNCS, vol. 7086. 119–134.
- FENNER, S. A., GASARCH, W. I., AND POSTOW, B. 2009. The Complexity of Finding SUBSEQ(A). *Theory of Computing Systems* 45, 3, 577–612.
- FILLIÁTRE, J.-C. 1997. Finite Automata Theory in Coq: A Constructive Proof of Kleene’s Theorem. Research Report 97–04, LIP - ENS Lyon.
- HAFTMANN, F. 2009. Code Generation from Specifications in Higher-Order Logic. Ph.D. thesis, Technical University of Munich.
- HAINES, L. H. 1969. On Free Monoids Partially Ordered by Embedding. *Journal of Combinatorial Theory* 6, 94–98.
- HARPER, R. 1999. Proof-Directed Debugging. *Journal of Functional Programming* 9, 4, 463–469.
- HOPCROFT, J. E. AND ULLMAN, J. D. 1969. *Formal Languages and Their Relation to Automata*. Addison-Wesley.
- KOZEN, D. 1997. *Automata and Computability*. Springer Verlag.
- KRAUSS, A. AND NIPKOW, T. 2012. Proof Pearl: Regular Expression Equivalence and Relation Algebra. *Journal of Automated Reasoning* 49, 1, 95–106.
- LAMMICH, P. AND TUERK, T. 2012. Applying Data Refinement for Monadic Programs to Hopcroft’s Algorithm. In *Proc. of the 3rd International Conference on Interactive Theorem Proving*. LNCS, vol. 7406. 166–182.
- NIPKOW, T. 1998. Verified Lexical Analysis. In *Proc. of the 11th International Conference on Theorem Proving in Higher Order Logics*. LNCS, vol. 1479. 1–15.
- NIPKOW, T. 2011. Gauss-Jordan Elimination for Matrices Represented as Functions. In *The Archive of Formal Proofs*, G. Klein, T. Nipkow, and L. Paulson, Eds. <http://afp.sourceforge.net/entries/Gauss-Jordan-Elim-Fun.shtml>. Formal proof development.
- OWENS, S., REPPY, J., AND TURON, A. 2009. Regular-Expression Derivatives Re-Examined. *Journal of Functional Programming* 19, 2, 173–190.
- OWENS, S. AND SLIND, K. 2008. Adapting Functional Programs to Higher Order Logic. *Higher-Order and Symbolic Computation* 21, 4, 377–409.
- ROSENBERG, A. L. 2006. A Big Ideas Approach to the Theory of Computation. Course notes for CMPSCI 401 at the University of Massachusetts.
- SAKAROVITCH, J. 2009. *Elements of Automata Theory*. Cambridge University Press.
- SHALLIT, J. 2008. *A Second Course in Formal Languages and Automata Theory*. Cambridge University Press.
- WU, C., ZHANG, X., AND URBAN, C. 2011a. A Formalisation of the Myhill-Nerode Theorem based on Regular Expressions (Proof Pearl). In *Proc. of the 2nd International Conference on Interactive Theorem Proving*. LNCS, vol. 6898. 341–356.
- WU, C., ZHANG, X., AND URBAN, C. 2011b. The Myhill-Nerode Theorem based on Regular Expressions. In *The Archive of Formal Proofs*, G. Klein, T. Nipkow, and L. Paulson, Eds. <http://afp.sourceforge.net/entries/Myhill-Nerode.shtml>. Formal proof development.
- YI, K. 2006. Educational Pearl: ‘Proof-Directed Debugging’ Revisited for a First-Order Version. *Journal of Functional Programming* 16, 6, 663–670.

## A. APPENDIX\*

PROOF OF LEMMA 2.1. For the right-to-left direction we assume  $X = B \cdot A^*$  and show that  $X = X \cdot A \cup B$  holds. From Property 2.1(i) we have  $A^* = A \cdot A^* \cup \{\epsilon\}$ , which is equal to  $A^* = A^* \cdot A \cup \{\epsilon\}$ . Adding  $B$  to both sides gives  $B \cdot A^* = B \cdot (A^* \cdot A \cup \{\epsilon\})$ , whose right-hand side is equal to  $(B \cdot A^*) \cdot A \cup B$ . Applying the assumed equation completes this direction.

For the other direction we assume  $X = X \cdot A \cup B$ . By a simple induction on  $n$ , we can establish the property

$$(*) \quad X = X \cdot A^{n+1} \cup \left( \bigcup_{m \leq n} B \cdot A^m \right)$$

Using this property we can show that  $B \cdot A^n \subseteq X$  holds for all  $n$ . From this we can infer  $B \cdot A^* \subseteq X$  using the definition of  $*$ . For the inclusion in the other direction we assume a string  $s$  with length  $k$  is an element in  $X$ . Since  $\epsilon \notin A$  we know by Property 2.1(ii) that  $s \notin X \cdot A^{k+1}$  since its length is only  $k$  (the strings in  $X \cdot A^{k+1}$  are all longer). From (\*) it follows then that  $s$  must be an element in  $\bigcup_{m \leq k} B \cdot A^m$ . This in turn implies that  $s$  is in  $\bigcup_n B \cdot A^n$ . Using Property 2.1(iii) this is equal to  $B \cdot A^*$ , as we needed to show.  $\square$

---

\*If the reviewers deem more suitable, the authors are prepared to drop material or move it to an electronic appendix.