

# A Formalisation of the Myhill-Nerode Theorem based on Regular Expressions (Proof Pearl)

Chunhan Wu<sup>1</sup>, Xingyuan Zhang<sup>1</sup>, and Christian Urban<sup>2</sup>

<sup>1</sup> PLA University of Science and Technology, China

<sup>2</sup> TU Munich, Germany

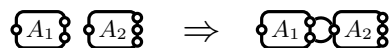
**Abstract.** There are numerous textbooks on regular languages. Nearly all of them introduce the subject by describing finite automata and only mentioning on the side a connection with regular expressions. Unfortunately, automata are difficult to formalise in HOL-based theorem provers. The reason is that they need to be represented as graphs, matrices or functions, none of which are inductive datatypes. Also convenient operations for disjoint unions of graphs and functions are not easily formalisable in HOL. In contrast, regular expressions can be defined conveniently as a datatype and a corresponding reasoning infrastructure comes for free. We show in this paper that a central result from formal language theory—the Myhill-Nerode theorem—can be recreated using only regular expressions.

## 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. [7]). It seems natural to exercise theorem provers by formalising the theorems and by verifying formally the algorithms.

There is however a problem: the typical approach to regular languages is to introduce finite automata and then define everything in terms of them. For example, a regular language is normally defined as one whose strings are recognised by a finite deterministic automaton. 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 HOL-based theorem prover, in our case Isabelle/HOL. Automata are built up from states and transitions that need to be represented as graphs, matrices or functions, none of which 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. Even worse, reasoning about graphs and matrices can be a real hassle in HOL-based theorem provers. 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 automata and hence will not be able to state certain properties about *all* automata, since there is no type quantification available in HOL (unlike in Coq, for example). An alternative, which provides us with a single type for 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 very adhoc constructions, which are not pleasant to reason about.

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 [9] dismisses for this the option of using identities, because it leads according to him to “messy proofs”. He 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.”*

Moreover, it is not so clear how to conveniently impose a finiteness condition upon functions in order to represent *finite* automata. The best is probably to resort to more advanced reasoning frameworks, such as *locales* or *type classes*, which are *not* available in all HOL-based theorem provers.

Because of these problems to do with representing automata, there seems to be no substantial formalisation of automata theory and regular languages carried out in HOL-based theorem provers. Nipkow [9] establishes the link between regular expressions and automata in the context of lexing. Berghofer and Reiter [2] formalise automata working over bit strings in the context of Presburger arithmetic. The only larger formalisations of automata theory are carried out in Nuprl [4] and in Coq [5].

In this paper, we will not attempt to formalise automata theory in Isabelle/HOL, but take a different approach to regular languages. Instead of defining a regular language as one where there exists an automaton that recognises all strings of the language, we define a regular language as:

**Definition 1.** *A language  $A$  is regular, provided there is a regular expression that matches all strings of  $A$ .*

The reason is that regular expressions, unlike graphs, matrices and functions, can be easily defined as inductive datatype. Consequently a corresponding reasoning infrastructure comes for free. This has recently been exploited in HOL4 with a formalisation of regular expression matching based on derivatives [11] and with an equivalence

checker for regular expressions in Isabelle/HOL [8]. The 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.

**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. We prove the part of this theorem stating that a regular expression has only finitely many partitions using certain tagging-functions. Again to our best knowledge, these tagging-functions have not been used before to establish 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  $[]$ . *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

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

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 1.

- (i)  $A^* = \{[]\} \cup A \cdot A^*$
- (ii) If  $[] \notin A$  and  $s \in A^{n+1}$  then  $n < |s|$ .
- (iii)  $B \cdot (\bigcup_n A^n) = (\bigcup_n B \cdot A^n)$

In (ii) we use the notation  $|s|$  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  $n$ . We omit the proofs for these properties, but invite the reader to consult our formalisation.<sup>3</sup>

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  $\{y \mid y \approx x\}$ .

Central to our proof will be the solution of equational systems involving equivalence classes of languages. For this we will use Arden’s Lemma [3], which solves equations of the form  $X = A \cdot X \cup B$  provided  $[] \notin A$ . However we will need the following ‘reverse’ version of Arden’s Lemma (‘reverse’ in the sense of changing the order of  $A \cdot X$  to  $X \cdot A$ ).

<sup>3</sup> Available at <http://www4.in.tum.de/~urbanc/regexp.html>

**Lemma 1 (Reverse Arden's Lemma).**

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

*Proof.* For the right-to-left direction we assume  $X = B \cdot A^*$  and show that  $X = X \cdot A \cup B$  holds. From Prop. 1(i) we have  $A^* = \{\square\} \cup A \cdot A^*$ , which is equal to  $A^* = \{\square\} \cup A^* \cdot A$ . Adding  $B$  to both sides gives  $B \cdot A^* = B \cdot (\{\square\} \cup A^* \cdot A)$ , whose right-hand side is equal to  $(B \cdot A^*) \cdot A \cup B$ . This 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

$$(*) \quad X = X \cdot A^{n+1} \cup \left( \bigcup_{m \in \{0..n\}} B \cdot A^m \right)$$

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  $\star$ . For the inclusion in the other direction we assume a string  $s$  with length  $k$  is an element in  $X$ . Since  $\square \notin A$  we know by Prop. 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  $\bigcup_{m \in \{0..k\}} B \cdot A^m$ . This in turn implies that  $s$  is in  $\bigcup_n B \cdot A^n$ . Using Prop. 1(iii) this is equal to  $B \cdot A^*$ , as we needed to show.  $\square$

Regular expressions are defined as the inductive datatype

$$r ::= \text{NULL} \mid \text{EMPTY} \mid \text{CHAR } c \mid \text{SEQ } r_1 r_2 \mid \text{ALT } r_1 r_2 \mid \text{STAR } r$$

and the language matched by a regular expression is defined as

$$\begin{aligned} \mathcal{L}(\text{NULL}) &\stackrel{\text{def}}{=} \emptyset & \mathcal{L}(\text{SEQ } r_1 r_2) &\stackrel{\text{def}}{=} \mathcal{L}(r_1) \cdot \mathcal{L}(r_2) \\ \mathcal{L}(\text{EMPTY}) &\stackrel{\text{def}}{=} \{\square\} & \mathcal{L}(\text{ALT } r_1 r_2) &\stackrel{\text{def}}{=} \mathcal{L}(r_1) \cup \mathcal{L}(r_2) \\ \mathcal{L}(\text{CHAR } c) &\stackrel{\text{def}}{=} \{[c]\} & \mathcal{L}(\text{STAR } r) &\stackrel{\text{def}}{=} \mathcal{L}(r)^* \end{aligned}$$

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  $rs$ . We only need to know the existence of such a regular expression and therefore we use Isabelle/HOL's *fold\_graph* and Hilbert's  $\varepsilon$  to define  $\dagger rs$ . This operation, roughly speaking, folds *ALT* over the set  $rs$  with *NULL* for the empty set. We can prove that for a finite set  $rs$

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

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

### 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 tertiary relation.

**Definition 2 (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

$$X_1 = \{\emptyset\} \quad X_2 = \{[c]\} \quad X_3 = UNIV - \{\emptyset, [c]\}$$

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 1.** *If finite  $(UNIV // \approx_A)$  then  $\exists r. A = \mathcal{L}(r)$ .*

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

$$\text{finals } A \stackrel{\text{def}}{=} \{[s]_{\approx_A} \mid s \in A\} \quad (3)$$

In our running example,  $X_2$  is the only equivalence class in *finals*  $\{[c]\}$ . It is straightforward to show that in general  $A = \bigcup \text{finals } A$  and  $\text{finals } A \subseteq UNIV // \approx_A$  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 Thm. 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

$$Y \stackrel{c}{\Longrightarrow} X \stackrel{\text{def}}{=} Y \cdot \{[c]\} \subseteq X \quad (4)$$

which means that if we concatenate 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  $X_1 \stackrel{c}{\Longrightarrow} X_2$ ,  $X_1 \stackrel{d}{\Longrightarrow} X_3$  with  $d$  being any other character than  $c$ , and  $X_3 \stackrel{d}{\Longrightarrow} X_3$  for any  $d$ .

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 equational system

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

where the terms  $(Y_{ij}, CHAR c_{ij})$  stand for all transitions  $Y_{ij} \xrightarrow{c_{ij}} X_i$ . There can only be finitely many terms of the form  $(Y_{ij}, CHAR c_{ij})$  in a right-hand side since by assumption there are only finitely many equivalence classes and only finitely many characters. The term  $\lambda(EMPTY)$  in the first equation acts as a marker for the initial state, that is the equivalence class containing  $\lambda$ .<sup>4</sup> 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  $X_{2..n}$  that the following equations

$$X_i = \mathcal{L}(Y_{i1}, CHAR c_{i1}) \cup \dots \cup \mathcal{L}(Y_{iq}, CHAR c_{iq}). \quad (5)$$

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

$$X_1 = \mathcal{L}(Y_{11}, CHAR c_{11}) \cup \dots \cup \mathcal{L}(Y_{1p}, CHAR c_{1p}) \cup \mathcal{L}(\lambda(EMPTY)). \quad (6)$$

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

$$\begin{aligned} Init\_rhs\ CS\ X &\stackrel{def}{=} \text{if } \lambda \in X \\ &\quad \text{then } \{(Y, CHAR\ c) \mid Y \in CS \wedge Y \xrightarrow{c} X\} \cup \{\lambda(EMPTY)\} \\ &\quad \text{else } \{(Y, CHAR\ c) \mid Y \in CS \wedge Y \xrightarrow{c} X\} \\ Init\ CS &\stackrel{def}{=} \{(X, Init\_rhs\ CS\ X) \mid X \in CS\} \end{aligned} \quad (7)$$

Because we use sets of terms for representing the right-hand sides of equations, we can prove (5) and (6) more concisely as

**Lemma 2.** *If  $(X, rhs) \in Init\ (UNIV // \approx_A)$  then  $X = \bigcup \mathcal{L}\ 'rhs$ .*

Our proof of Thm. 1 will proceed by transforming the initial equational system into one in *solved form* maintaining the invariant in Lem. 2. 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

<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 [3], 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 reverse version of Arden's Lemma.

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

$$(Y, r_2) \triangleleft r_1 \stackrel{def}{=} (Y, SEQ\ r_2\ r_1) \quad \lambda(r_2) \triangleleft r_1 \stackrel{def}{=} \lambda(SEQ\ r_2\ r_1)$$

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:

$$\begin{aligned} Arden\ X\ rhs &\stackrel{def}{=} \text{let} \\ &\quad rhs' = rhs - \{(X, r) \mid (X, r) \in rhs\} \\ &\quad r' = STAR(\dagger\{r \mid (X, r) \in rhs\}) \\ &\text{in } rhs' \triangleleft r' \end{aligned} \tag{8}$$

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  $rhs'$ . It can be easily seen that this 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

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

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

**Lemma 3.** *Given an equation  $X = rhs$ . If  $X = \bigcup \mathcal{L}'\ rhs$ , ardenable  $rhs$ , and finite  $rhs$ , then  $X = \bigcup \mathcal{L}'\ (Arden\ X\ rhs)$ .*

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  $rhs$ .

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

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  $X$ .

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.

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

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 [1]), 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 (7) using the principle:

$$\begin{array}{l} \text{invariant } (\text{Init } (UNIV // \approx_A)) \\ \forall ES. \text{invariant } ES \wedge \text{Cond } ES \longrightarrow \text{invariant } (\text{Iter } X \ ES) \\ \forall ES. \text{invariant } ES \wedge \text{Cond } ES \longrightarrow \text{card } (\text{Iter } X \ ES) < \text{card } ES \\ \forall ES. \text{invariant } ES \wedge \neg \text{Cond } ES \longrightarrow P \ ES \\ \hline P (\text{Solve } X (\text{Init } (UNIV // \approx_A))) \end{array} \quad (9)$$

This principle states that given an invariant (which we will specify below) we can prove a property  $P$  involving *Solve*. For this we have to discharge the following proof obligations: first the initial equational system satisfies the invariant; second the iteration step *Iter* preserves the invariant as long as the condition *Cond* holds; third *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  $\text{Solve } X (\text{Init } (UNIV // \approx_A))$  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:



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

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 4.** *If finite  $(UNIV // \approx_A)$  then invariant  $(Init (UNIV // \approx_A))$ .*

*Proof.* Finiteness is given by the assumption and the way how we set up the initial equational system. Soundness is proved in Lem. 2. 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 *Iter* preserves the invariant.

**Lemma 5.** *If invariant  $ES$  and  $(X, rhs) \in ES$  and  $Cond \text{ } ES$  then invariant  $(Iter \text{ } X \text{ } ES)$ .*

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

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

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 Lem. 3). 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 6.** *If invariant  $ES$  and  $(X, rhs) \in ES$  and  $Cond\ ES$  then  $card\ (Iter\ X\ ES) < card\ ES$ .*

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

*Proof.* In order to prove this lemma using (9), we have to use a slightly stronger invariant since Lem. 5 and 6 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 Lem. 4, the more general invariant holds for the initial equational system. This is premise 1 of (9). Premise 2 is given by Lem. 5 and the fact that *Iter* might modify the  $rhs$  in the equation  $X = rhs$ , but does not remove it. Premise 3 of (9) is by Lem. 6. 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  $Solve\ X\ (Init\ (UNIV // \approx_A)) = \{(X, rhs)\}$  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 8.** *If finite  $(UNIV // \approx_A)$  and  $X \in UNIV // \approx_A$  then  $\exists r. X = \mathcal{L}(r)$ .*

*Proof.* By the preceding lemma, we know that there exists a  $rhs$  such that  $Solve\ X\ (Init\ (UNIV // \approx_A))$  returns the equation  $X = rhs$ , and that the invariant holds for this equation. That means we know  $X = \bigcup \mathcal{L}\ 'rhs$ . We further know that this is equal to  $\bigcup \mathcal{L}\ '(Arden\ X\ rhs)$  using the properties of the invariant and Lem. 3. Using the validity property for the equation  $X = rhs$ , we can infer that  $rhs\ rhs \subseteq \{X\}$  and because the *Arden* operation removes that  $X$  from  $rhs$ , that  $rhss\ (Arden\ X\ rhs) = \emptyset$ . 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  $X = \mathcal{L}(+rs)$ . With this we can conclude the proof.  $\square$

Lem. 8 allows us to finally give a proof for the first direction of the Myhill-Nerode theorem.

*Proof (of Thm. 1).* By Lem. 8 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  $\bigcup finals A = \mathcal{L}(\dagger rs)$ . Since the left-hand side is equal to  $A$ , we can use  $\dagger rs$  as the regular expression that is needed in the theorem.  $\square$

## 4 Myhill-Nerode, Second Part

We will prove in this section the second part of the Myhill-Nerode theorem. It can be formulated in our setting as follows:

**Theorem 2.** *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 *NULL*, *EMPTY* and *CHAR* are routine, because we can easily establish that

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

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 solve directly. The reader is invited to try.

Our proof 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  $range f \stackrel{def}{=} f ' UNIV$ ). With this we will be able to infer that the tagging-functions, seen as relations, give rise to finitely many equivalence classes of  $UNIV$ . 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 (a relation  $R_1$  is said to *refine*  $R_2$  provided  $R_1 \subseteq R_2$ ). We formally define the notion of a *tagging-relation* as follows.

**Definition 3 (Tagging-Relation).** *Given a tagging-function  $tag$ , then two strings  $x$  and  $y$  are tag-related provided*

$$x =_{tag} y \stackrel{def}{=} tag x = tag y .$$

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

$$\text{If finite } (f ' A) \text{ and inj\_on } f A \text{ then finite } A. \quad (10)$$

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 9.** *If finite (range tag) then finite (UNIV // =tag=).*

*Proof.* We set in (10),  $f$  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 // =tag=) is a subset of  $\text{range } f$ , and so also finite. Injectivity amounts to showing that  $X = Y$  under the assumptions that  $X, Y \in \text{UNIV // =tag=}$  and  $f X = f Y$ . From the assumptions we can 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.  $\square$

**Lemma 10.** *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 (10). This time we set  $f$  to be  $X \mapsto \{\llbracket x \rrbracket_{R_1} \mid x \in X\}$ . It is easy to see that  $\text{finite}(f \text{ UNIV // } R_2)$  because it is a subset of  $\text{Pow}(\text{UNIV // } R_2)$ , which is 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  $f X = f Y$ . 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  $\llbracket x \rrbracket_{R_1} \in f X$  and further  $\llbracket x \rrbracket_{R_1} \in f Y$ . 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 Lem. 9 and 10 together, means in order to show that  $\text{UNIV // } \approx_{\mathcal{L}(r)}$  is finite, we have to find a tagging-function whose range can be shown to be finite and whose tagging-relation refines  $\approx_{\mathcal{L}(r)}$ . Let us attempt the *ALT*-case first.

*Proof (ALT-Case).* We take as tagging-function

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

where  $A$  and  $B$  are some arbitrary languages. 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  $\text{tag}_{ALT} A B$  is a subset of this product set—so finite. It remains to be shown that  $=\text{tag}_{ALT} A B =$  refines  $\approx_{A \cup B}$ . This amounts to showing

$$\text{tag}_{ALT} A B x = \text{tag}_{ALT} A B y \longrightarrow x \approx_{A \cup B} y$$

which by unfolding the Myhill-Nerode relation is identical to

$$\forall z. \text{tag}_{ALT} A B x = \text{tag}_{ALT} A B y \wedge x @ z \in A \cup B \longrightarrow y @ z \in A \cup B \quad (11)$$

since both  $=\text{tag}_{ALT} A B =$  and  $\approx_{A \cup B}$  are symmetric. To solve (11) we just have to unfold the definition of the tagging-function and analyse in which set,  $A$  or  $B$ , the string  $x @ z$  is. The definition of the tagging-function will give us in each case the information to infer that  $y @ z \in A \cup B$ . Finally we can discharge this case by setting  $A$  to  $\mathcal{L}(r_1)$  and  $B$  to  $\mathcal{L}(r_2)$ .  $\square$

The pattern in (11) is repeated for the other two cases. Unfortunately, they are slightly more complicated. In the *SEQ*-case we essentially have to be able to infer that

$$\dots x @ z \in A \cdot B \longrightarrow y @ z \in A \cdot B$$

using the information given by the appropriate tagging-function. The complication is to find out what the possible splits of  $x @ z$  are to be in  $A \cdot B$  (this was easy in case of  $A \cup B$ ). To deal with this complication we define the notions of *string prefixes*

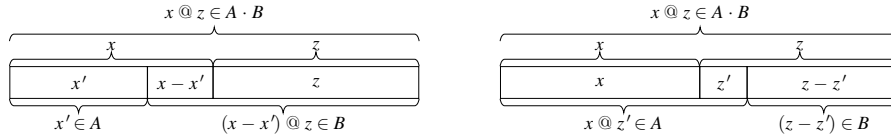
$$x \leq y \stackrel{\text{def}}{=} \exists z. y = x @ z \quad x < y \stackrel{\text{def}}{=} x \leq y \wedge x \neq y$$

and *string subtraction*:

$$\square - y \stackrel{\text{def}}{=} \square \quad x - \square \stackrel{\text{def}}{=} x \quad cx - dy \stackrel{\text{def}}{=} \text{if } c = d \text{ then } x - y \text{ else } cx$$

where  $c$  and  $d$  are characters, and  $x$  and  $y$  are strings.

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 there is a prefix of  $x$  in  $A$  and the rest is in  $B$  (first picture), or  $x$  and a prefix of  $z$  is in  $A$  and the rest in  $B$  (second picture). In both cases we have to show that  $y @ z \in A \cdot B$ . For this we use the following tagging-function

$$\text{tag}_{SEQ} A B x \stackrel{\text{def}}{=} (\llbracket x \rrbracket \approx_A, \{ \llbracket (x - x') \rrbracket \approx_B \mid x' \leq x \wedge x' \in A \})$$

with the idea that in the first split we have to make sure that  $(x - x') @ z$  is in the language  $B$ .

*Proof (SEQ-Case).* If *finite*  $(UNIV // \approx_A)$  and *finite*  $(UNIV // \approx_B)$  then *finite*  $(UNIV // \approx_A \times \text{Pow}(UNIV // \approx_B))$  holds. The range of  $\text{tag}_{SEQ} A B$  is a subset of this product set, and therefore finite. We have to show injectivity of this tagging-function as

$$\forall z. \text{tag}_{SEQ} A B x = \text{tag}_{SEQ} A B y \wedge x @ z \in A \cdot B \longrightarrow y @ z \in A \cdot B$$

There are two cases to be considered (see pictures above). First, there exists a  $x'$  such that  $x' \in A$ ,  $x' \leq x$  and  $(x - x') @ z \in B$  hold. We therefore have

$$\llbracket (x - x') \rrbracket \approx_B \in \{ \llbracket (x - x') \rrbracket \approx_B \mid x' \leq x \wedge x' \in A \}$$

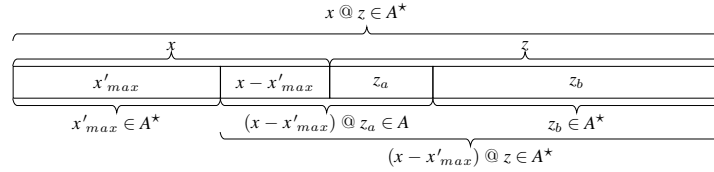
and by the assumption about  $\text{tag}_{SEQ} A B$  also

$$\llbracket (x - x') \rrbracket \approx_B \in \{ \llbracket (y - y') \rrbracket \approx_B \mid y' \leq y \wedge y' \in A \}$$

That means there must be a  $y'$  such that  $y' \in A$  and  $\llbracket (x - x') \rrbracket \approx_B = \llbracket (y - y') \rrbracket \approx_B$ . This equality means that  $x - x' \approx_B y - y'$  holds. Unfolding the Myhill-Nerode relation and together with the fact that  $(x - x') @ z \in B$ , we have  $(y - y') @ z \in B$ . We already know  $y' \in A$ , therefore  $y @ z \in A \cdot B$ , as needed in this case.

Second, there exists a  $z'$  such that  $x @ z' \in A$  and  $z - z' \in B$ . By the assumption about  $\text{tag}_{SEQ} A B$  we have  $\llbracket x \rrbracket \approx_A = \llbracket y \rrbracket \approx_A$  and thus  $x \approx_A y$ . Which means by the Myhill-Nerode relation that  $y @ z' \in A$  holds. Using  $z - z' \in B$ , we can conclude also in this case with  $y @ z \in A \cdot B$ . We again can complete the *SEQ*-case by setting  $A$  to  $\mathcal{L}(r_1)$  and  $B$  to  $\mathcal{L}(r_2)$ .  $\square$

The case for *STAR* is similar to *SEQ*, but poses a few extra challenges. When we analyse the case that  $x @ z$  is an element in  $A^*$  and  $x$  is not the empty string, we have the following picture:



We can find a strict prefix  $x'$  of  $x$  such that  $x' \in A^*$ ,  $x' < x$  and the rest  $(x - x') @ z \in A^*$ . For example the empty string  $\epsilon$  would do. 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  $x'_{max}$ . Now for the rest of the string  $(x - x'_{max}) @ z$  we know it is in  $A^*$ . By definition of  $A^*$ , we can separate this string into two parts, say  $a$  and  $b$ , such that  $a \in A$  and  $b \in A^*$ . Now  $a$  must be strictly longer than  $x - x'_{max}$ , otherwise  $x'_{max}$  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  $(x - x'_{max}) @ z_a \in A$  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  $(x - x'_{max}) @ z_a$ .

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

$$\text{tag}_{STAR} A x \stackrel{\text{def}}{=} \{[(x - x')] \approx_A \mid x' < x \wedge x' \in A^*\}$$

*Proof (STAR-Case).* If *finite* ( $UNIV // \approx_A$ ) then *finite* ( $Pow (UNIV // \approx_A)$ ) holds. The range of  $\text{tag}_{STAR} A$  is a subset of this set, and therefore finite. Again we have to show injectivity of this tagging-function as

$$\forall z. \text{tag}_{STAR} A x = \text{tag}_{STAR} A y \wedge x @ z \in A^* \longrightarrow y @ z \in A^*$$

We first need to consider the case that  $x$  is the empty string. From the assumption we can infer  $y$  is the empty string and 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 we have

$$[(x - x'_{max})] \approx_A \in \{[(x - x')] \approx_A \mid x' < x \wedge x' \in A^*\}$$

which by assumption is equal to

$$[(x - x'_{max})] \approx_A \in \{[(y - y')] \approx_A \mid y' < y \wedge y' \in A^*\}$$

and we know that we have a  $y' \in A^*$  and  $y' < y$  and also know  $x - x'_{max} \approx_A y - y'$ . Unfolding the Myhill-Nerode relation we know  $(y - y') @ z_a \in A$ . We also know that  $z_b \in A^*$ . Therefore  $y' @ ((y - y') @ z_a) @ z_b \in A^*$ , which means  $y @ z \in A^*$ . As the last step we have to set  $A$  to  $\mathcal{L}(r)$  and complete the proof.  $\square$

## 5 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 HOL-based 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 3 (The 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 [7]). 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 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. Proving the existence of 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.

While regular expressions are convenient in formalisations, 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 [10]. For an implementation of a simple regular expression matcher, whose correctness has been formally established, we refer the reader to Owens and Slind [11].

Our formalisation consists of 780 lines of Isabelle/Isar code for the first direction and 460 for the second, plus around 300 lines of standard material about regular languages. While this might be seen as too large to count as a concise proof pearl, this should be seen in the context of the work done by Constable et al [4] who formalised the Myhill-Nerode theorem in Nuprl using automata. They write that their four-member team needed something on the magnitude of 18 months for their formalisation. 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, who were able to follow the proofs from [6], we had to find our own arguments. So for us the formalisation was not the bottleneck. 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. The code of ours can be found in the Mercurial Repository at <http://www4.in.tum.de/~urbanc/regexp.html>.

Our proof of the first direction is very much inspired by *Brzozowski's algebraic method* used to convert a finite automaton to a regular expression [3]. The close con-

nection 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. Since we identify equivalence classes with the states of the automaton, then 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 briefly considered using the method Brzozowski presented in the Appendix of [3] in order to prove the second direction of the Myhill-Nerode theorem. There he calculates the derivatives for regular expressions and shows that for every language there can be only finitely many of them (if regarded equal modulo ACI). We could have used as tagging-function the set of derivatives of a regular expression with respect to a language. Using the fact that two strings are Myhill-Nerode related whenever their derivative is the same, together with the fact that there are only finitely such derivatives would give us a similar argument as ours. However it seems not so easy to calculate the set of derivatives modulo ACI. Therefore we preferred our direct method of using tagging-functions. This is also where our method shines, because we can completely side-step the standard argument [7] where automata need to be composed, which as stated in the Introduction is not so easy to formalise in a HOL-based theorem prover. However, it is also the direction where we had to spend most of the ‘conceptual’ time, as our proof-argument based on tagging-functions is new for establishing the Myhill-Nerode theorem. All standard proofs of this direction use arguments over automata.

## References

1. S. Berghofer and T. Nipkow. Executing Higher Order Logic. In *Proc. of the International Workshop on Types for Proofs and Programs*, volume 2277 of *LNCS*, pages 24–40, 2002.
2. S. Berghofer and M. Reiter. Formalizing the Logic-Automaton Connection. In *Proc. of the 22nd International Conference on Theorem Proving in Higher Order Logics*, volume 5674 of *LNCS*, pages 147–163, 2009.
3. J. A. Brzozowski. Derivatives of Regular Expressions. *J. ACM*, 11:481–494, 1964.
4. R. L. Constable, P. B. Jackson, P. Naumov, and J. C. Uribe. Constructively Formalizing Automata Theory. In *Proof, Language, and Interaction*, pages 213–238. MIT Press, 2000.
5. J.-C. Filliâtre. Finite Automata Theory in Coq: A Constructive Proof of Kleene’s Theorem. Research Report 97–04, LIP - ENS Lyon, 1997.
6. J. E. Hopcroft and J. D. Ullman. *Formal Languages and Their Relation to Automata*. Addison-Wesley, 1969.
7. D. Kozen. *Automata and Computability*. Springer Verlag, 1997.
8. A. Kraus and T. Nipkow. Proof Pearl: Regular Expression Equivalence and Relation Algebra. To appear in *Journal of Automated Reasoning*, 2011.
9. T. Nipkow. Verified Lexical Analysis. In *Proc. of the 11th International Conference on Theorem Proving in Higher Order Logics*, volume 1479 of *LNCS*, pages 1–15, 1998.
10. S. Owens, J. Reppy, and A. Turon. Regular-Expression Derivatives Re-Examined. *Journal of Functional Programming*, 19(2):173–190, 2009.
11. S. Owens and K. Slind. Adapting Functional Programs to Higher Order Logic. *Higher-Order and Symbolic Computation*, 21(4):377–409, 2008.