

Constructive Formalization of Regular Languages

Jan-Oliver Kaiser

September 13, 2012

Abstract

Existing formalizations of regular languages in constructive settings are mostly limited to regular expressions and finite automata. Furthermore, these usually require in the order of 10,000 lines of code. The goal of this thesis is to show that an extensive, yet elegant formalization of regular languages can be achieved in constructive type theory. In addition to regular expressions and finite automata, our formalization includes the Myhill-Nerode theorem. The entire development weighs in at approximately 3,300 lines of code.

Citations?

Reduce
& up-
date

Contents

1	Introduction	4
1.1	Recent work	4
2	Coq and SSReflect	5
2.1	CoQ	5
2.2	SSREFLECT	5
2.2.1	Finite Types and Ordinals	5
2.2.2	Boolean Reflection	6
2.2.3	Boolean Predicates	6
3	Decidable Languages	7
3.1	Definitions	7
3.1.1	Operations on languages	8
3.2	Regular Languages	10
3.2.1	Regular Expressions	10
4	Finite Automata	12
4.1	Definition	12
4.1.1	Non-Deterministic Finite Automata	12
4.1.2	Deterministic Finite Automata	13
4.2	Connected Components	16
4.3	Emptiness	18
4.4	Deciding Equivalence of Finite Automata	18
4.5	Regular Expressions and Finite Automata	19
4.5.1	Regular Expressions to Finite Automata	19
4.5.2	Deciding Equivalence of Regular Expressions	19
4.5.3	Finite Automata to Regular Expressions	19
5	Myhill-Nerode	23
5.1	Definition	23
5.2	Finite Partitionings and Equivalence Classes	24
5.3	Minimizing Equivalence Classes	25
5.4	Finite Automata and Myhill-Nerode	26
5.4.1	Finite Automata to Myhill-Nerode	26

<i>CONTENTS</i>	3
5.4.2 Myhill-Nerode to Finite Automata	26
6 Conclusion	27

Chapter 1

Introduction

Regular languages are a well-studied class of formal languages. We will prove the equivalence of three well-known characterizations of regular languages: regular expressions, finite automata and the characterization given by Myhill-Nerode theorem.

History

Theoretical
importance

Practical
importance?

1.1 Recent work

There have been many publications on regular languages in recent years. Many of them investigate decidability of equivalence of regular languages, though there have also been new equivalence proofs regarding different characterizations of regular languages.

Chapter 2

Coq and SSReflect

We decided to employ the Small Scale Reflection Extension (**SSReflect**¹) for the **Coq**² proof assistant. The most important factors in this decision were SSREFLECT's excellent support for finite types, list operations and graphs. SSREFLECT also introduces an alternative scripting language that can often be used to shorten the bookkeeping overhead of proofs considerably.

2.1 Coq

Description,
citation

2.2 SSReflect

SSREFLECT is a set of extensions to the proof scripting language of the COQ proof assistant. They were originally developed to support small-scale reflection. However, most of them are of quite general nature and improve the functionality of COQ in most basic areas such as script layout and structuring, proof context management and rewriting [4].

2.2.1 Finite Types and Ordinals

The most important feature of SSREFLECT for our purpose are finite types. SSREFLECT provides boolean versions of the universal and existential quantifiers on finite types, **forallb** and **existsb**. We can compute the number of elements in a finite type F with $\#|F|$. **enum** gives a list of all items of a finite type. Finite types also come with enumeration functions which provide a consistent ordering. The corresponding functions are **enum.rank** and

¹<http://www.msr-inria.inria.fr/Projects/math-components>

²<http://coq.inria.fr/>

`enum_val`. The input of `enum_val` and the result of `enum_rank` are ordinals, i.e. values in $[0, \#|F|-1]$. The corresponding type can be written as $1_{\#|F|}$.

2.2.2 Boolean Reflection

SSREFLECT offers boolean reflections for decidable propositions. This allows us to switch back and forth between equivalent boolean and propositional predicates.

2.2.3 Boolean Predicates

We make use of SSREFLECT's syntax to specify boolean predicates. This allows us to specify predicates in a way that resembles set-theoretic notation, e.g. `[pred x | <boolean expression in x>]`. Furthermore, we can use the functions `pred1` and `pred0` to specify the singleton predicate and the empty predicate, respectively. The complement of a predicate can be written as `[predC p]`. The syntax for combining predicates is `[pred? p1 & p2]`, with `?` being replaced with one of `U` (union), `I` (intersection) or `D` (difference). There is also syntax for the preimage of a predicate under a function which can be written as `[preim f of p]`.

There are also applicative (functional) versions of `predC`, `predU`, `predI`, `predD` which are functions that take predicates and return predicates.

Chapter 3

Decidable Languages

We give basic definitions for languages, operators on languages and, finally, regular languages. We provide the corresponding formalizations from our development and prove their correctness.

3.1 Definitions

We closely follow the definitions from [5]. An **alphabet** Σ is a finite set of symbols. A **word** w is a finite sequence of symbols chosen from some alphabet. We use $|w|$ to denote the **length** of a word w . ε denotes the empty word. Given two words $w_1 = a_1 \cdots a_n$ and $w_2 = b_1 \cdots b_m$, the **concatenation** of w_1 and w_2 is defined as $a_1 \cdots a_n \cdot b_1 \cdots b_m$ and denoted $w_1 \cdot w_2$ or just $w_1 w_2$. A **language** is a set of words. The **residual language** of a language L with respect to symbol a is the set of words u such that au is in L . The residual is denoted $res_a(L)$. We define Σ^k to be the **set of words of length k**. The **set of all words** over an alphabet Σ is denoted Σ^* , i.e., $\Sigma^* = \bigcup_{k \in \mathbb{N}} \Sigma^k$. A language L is **decidable** if and only if there exists a boolean predicate that decides membership in L . We will only deal with decidable languages from here on. Throughout the remaining document, we will assume a fixed alphabet Σ .

We employ finite types to formalize alphabets. In the most definitions, alphabets will not be made explicit. However, the same name and type will be used throughout the entire development. Words are formalized as sequences over the alphabet. Decidable languages are represented by functions from *word* to *bool*.

Variable `char`: `finType`.

Definition `word` := `seq char`.

Definition `language` := `pred word`.

Definition `residual x L : language` := `[preim cons x of L]`.

3.1.1 Operations on languages

The **concatenation** of two languages L_1 and L_2 is denoted $L_1 \cdot L_2$ and is defined as the set of words $w = w_1 w_2$ such that w_1 is in L_1 and w_2 is in L_2 . The **Kleene closure** of a language L (also called Kleene star) is denoted L^* and is defined as the set of words $w = w_1 \cdots w_k$ such that $w_1 \dots w_k$ are in L . ε is contained in L^* ($k = 0$). We define the **complement** of a language L as $L \setminus \Sigma^*$, which we write as $\neg L$. Furthermore, we make use of the standard set operations **union** and **intersection**.

For our COQ development, take Coquand and Siles's [2] implementation of these operators. `plus` and `prod` refer to union and intersection, respectively. Additionally, we also introduce the singleton languages (`atom`), the empty language (`void`) and the language containing only the empty word (`eps`).

Definition `conc L1 L2 : language :=`

`fun v => existsb i : 'L (size v).+1, L1 (take i v) && L2 (drop i v).`

Definition `star L : language :=`

`fix star v := if v is x :: v' then conc (residual x L) star v' else true.`

Definition `compl L : language := predC L.`

Definition `plus L1 L2 : language := [predU L1 & L2].`

Definition `prod L1 L2 : language := [predI L1 & L2].`

Definition `atom x : language := pred1 [:: x].`

Definition `void : language := pred0.`

Definition `eps : language := pred1 [::].`

The definition of `conc` is based on a characteristic property of the concatenation of two languages. The following lemma proves this property.

Lemma 3.1.1. *Let $L_1, L_2, w = a_1 \cdots a_k$ be given. We have that*

$$w \in L_1 \cdot L_2 \iff \exists n \in \mathbb{N}. 0 < n \leq k \wedge a_1 \cdots a_{n-1} \in L_1 \wedge a_n \cdots a_k \in L_2.$$

Proof. “ \Rightarrow ” From $w \in L_1 \cdot L_2$ we have w_1, w_2 such that $w = w_1 w_2 \wedge w_1 \in L_1 \wedge w_2 \in L_2$. We choose $n := |w_1| + 1$. We then have that $a_1 \cdots a_{n-1} = a_1 \cdots a_{|w_1|} = w_1$ and $w_1 \in L_1$ by assumption. Similarly, $a_n \cdots a_k = a_{|w_1|+1} \cdots a_k = w_2$ and $w_2 \in L_2$ by assumption.

“ \Leftarrow ” We choose $w_1 := a_1 \cdots a_{n-1}$ and $w_2 := a_n \cdots a_k$. By assumption we have that $w = w_1 w_2$. We also have that $a_1 \cdots a_{n-1} \in L_1$ and $a_n \cdots a_k \in L_2$. It follows that $w_1 \in L_1$ and $w_2 \in L_2$. \square

Listing 3.1: Formalization of lemma 3.1.1

Lemma `concP : forall {L1 L2 v},`

`reflect (exists2 v1, v1 \in L1 & exists2 v2, v2 \in L2 & v = v1 ++ v2)`
`(v \in conc L1 L2).`

The implementation of `star` makes use of a characteristic property of the Kleene closure, which is that any nonempty word in L^* can be seen as the concatenation of a nonempty word in L and a (possibly empty) word in L^* . The following lemma proves the correctness of this property.

Lemma 3.1.2. *Let $L, w = a_1 \cdots a_k$ be given. We have that*

$$w \in L^* \iff \begin{cases} a_2 \cdots a_k \in \text{res}_{a_1}(L) \cdot L^*, & \text{if } |w| > 0; \\ w = \varepsilon, & \text{otherwise.} \end{cases}$$

Proof. “ \Rightarrow ” We do a case distinction on $|w| = 0$.

1. $|w| = 0$. It follows that $w = \varepsilon$.
2. $|w| \neq 0$, i.e. $|w| > 0$. From $w \in L^*$ we have $w = w_1 \cdots w_l$ such that $w_1 \cdots w_l$ are in L . There exists a minimal n such that $|w_n| > 0$ and for all $m < n$, $|w_m| = 0$. Let $w_n = b_1 \cdots b_p$. We have that $b_2 \cdots b_p \in \text{res}_{b_1}(L)$. Furthermore, we have that $w_{n+1} \cdots w_l \in L^*$. We also have $a_1 = b_1$ and $w = a_1 \cdots a_k = w_n \cdots w_l$. Therefore, we have $a_2 \cdots a_k \in \text{res}_{a_1}(L) \cdot L^*$.

“ \Leftarrow ” We do a case distinction on the disjunction.

1. $w = \varepsilon$. Then $w \in L^*$ by definition.
2. $a_2 \cdots a_k \in \text{res}_{a_1}(L) \cdot L^*$. By lemma 3.1.1 we have n such that $a_2 \cdots a_{n-1} \in \text{res}_{a_1}(L)$ and $a_n \cdots a_k \in L^*$. By definition of `res`, we have $a_1 \cdots a_{n-1} \in L$. Furthermore, we also have $a_n \cdots a_k = w_1 \cdots w_l$ such that $w_1 \cdots w_l$ are in L . We choose $w_0 := a_1 \cdots a_{n-1}$. It follows that $w = w_0 w_1 \cdots w_l$ with w_0, w_1, \dots, w_l in L . Therefore, $w \in L^*$.

□

The formalization of lemma 3.1.2 connects the characteristic property used to define `star` to the mathematical definition. The propositional formula given here appears slightly more restrictive than our mathematical definition as it requires all words from L to be nonempty. Mathematically, however, this is no restriction. The advantage of this formalization is can be seen in later proofs, where induction over the length of a word $w \in L^*$ is necessary. In such places, we will use lemma 3.1.2 to get $w_1 \dots w_k \in L$ such that $w = w_1 \cdots w_k$ and proceed by induction over the list of words $w_1 \dots w_k$. Since every word w_i has a length of at least one, we shorten w by at least one character. Thus, we can use the inductive hypothesis of the induction over $|w|$ to prove the claim.

Listing 3.2: Formalization of lemma 3.1.2

```
Lemma starP : forall {L v},
  reflect (exists2 vv, all [predD L & eps] vv & v = flatten vv)
    (v \in star L).
```

Theorem 3.1.1. *The decidable languages are closed under concatenation, Kleene star, union, intersection and complement.*

Proof. We have already given algorithms for every operator. It remains to show that they are correct. For concatenation and the Kleene star, we have shown in lemma 3.1.1 and lemma 3.1.2 that the formalizations are equivalent to the mathematical definitions. The remaining operators (union, intersection, complement) can be applied directly to the result of the languages' boolean decision functions. \square

3.2 Regular Languages

Definition 3.2.1. *The set of regular languages REG is defined to be exactly those languages generated by the following inductive definition:*

$$\begin{array}{c} \overline{\emptyset \in REG} \qquad \overline{\{\varepsilon\} \in REG} \qquad \frac{a \in \Sigma}{\{a\} \in REG} \qquad \frac{L \in REG}{L^* \in REG} \\[10pt] \frac{L_1 \in REG \quad L_2 \in REG}{L_1 \cup L_2 \in REG} \qquad \frac{L_1 \in REG \quad L_2 \in REG}{L_1 \cdot L_2 \in REG} \end{array}$$

3.2.1 Regular Expressions

Regular expressions mirror the definition of regular languages very closely. We will consider **extended regular expressions** that include negation (Not), intersection (And) and a single-symbol wildcard (Dot). Therefore, we have the following syntax for regular expressions:

$$r, s := \emptyset \mid \varepsilon \mid \cdot \mid a \mid r^* \mid r + s \mid r \& s \mid rs \mid \neg r$$

The semantics of these constructors are as follows:

1. $\mathcal{L}(\emptyset) = \emptyset$
2. $\mathcal{L}(\varepsilon) = \{\varepsilon\}$
3. $\mathcal{L}(\cdot) = \Sigma$
4. $\mathcal{L}(a) = \{a\}$
5. $\mathcal{L}(r^*) = \mathcal{L}(r)^*$
6. $\mathcal{L}(r + s) = \mathcal{L}(r) \cup \mathcal{L}(s)$

$$7. \mathcal{L}(r \& s) = \mathcal{L}(r) \cap \mathcal{L}(s)$$

$$8. \mathcal{L}(rs) = \mathcal{L}(r) \cdot \mathcal{L}(s)$$

Definition 3.2.2. *We say that two regular expressions r and s are equivalent if and only if*

$$\mathcal{L}(r) = \mathcal{L}(s).$$

We will later show that equivalence of regular expressions is decidable.

We take the implementation of regular expressions from Coquand and Siles's development ([2]), which is also based on SSREFLECT and comes with helpful infrastructure for our proofs. They also give a decision function the language associated with a regular expression.

Listing 3.3: Regular Expressions

```

Inductive regular_expression :=
| Void
| Eps
| Dot
| Atom of symbol
| Star of regular_expression
| Plus of regular_expression & regular_expression
| And of regular_expression & regular_expression
| Conc of regular_expression & regular_expression
| Not of regular_expression .

Fixpoint mem_reg e :=
  match e with
  | Void => void
  | Eps => eps
  | Dot => dot
  | Atom x => atom x
  | Star e1 => star (mem_reg e1)
  | Plus e1 e2 => plus (mem_reg e1) (mem_reg e2)
  | And e1 e2 => prod (mem_reg e1) (mem_reg e2)
  | Conc e1 e2 => conc (mem_reg e1) (mem_reg e2)
  | Not e1 => compl (mem_reg e1)
  end.

```

We will later prove that this definition is equivalent to the inductive definition of regular languages in 3.2.1. In order to do that, we introduce a predicate on regular expressions that distinguishes **standard regular expressions** from **extended regular expressions** (as introduced above). The grammar of standard regular expression is as follows:

$$r, s := \emptyset \mid \varepsilon \mid a \mid r^* \mid r + s \mid rs$$

Connect
stan-
dard reg-
exp to
reg. lan-
guages

Chapter 4

Finite Automata

Another way of characterizing regular languages are finite automata (FA). We will show that the languages of finite automata are exactly the regular languages. Furthermore, we will also derive a decision procedure for equivalence of regular expressions.

4.1 Definition

A finite automaton consists of

1. finite set of states Q ,
2. an alphabet Σ ,
3. a starting state $s \in Q$,
4. a set of final states $F \subseteq Q$
5. and a state-transition relation δ . [5]

We define a **run** of a word $w \in \Sigma^*$ on an automaton $A = (\Sigma, Q, s, F, \delta)$ as any sequence of states σ such that $\forall 0 \leq i < |\sigma| - 1. (\sigma_i, w_i, \sigma_{i+1}) \in \delta$. A word w is **accepted** by A if and only if there exists a run σ of w on A such that $\sigma_0 = s \wedge \sigma_{|\sigma|-1} \in F$. The **language** of A is exactly the set of words accepted by A and is denoted $\mathcal{L}(A)$. It will later be useful to also have an acceptance criterion defined by runs starting in a given state $x \in Q$, for which we will denote the resulting language $\mathcal{L}_x(A)$.

4.1.1 Non-Deterministic Finite Automata

Finite automata can be **non-deterministic** (NFA) in the sense that there exist multiple distinct runs for a word. This is the case if and only if δ is not functional.

Listing 4.1: Non-Deterministic Finite Automata

```

Record nfa : Type :=
{
  nfa_state :> finType;
  nfa_s : nfa_state;
  nfa_fin : pred nfa_state;
  nfa_step : nfa_state -> char -> pred nfa_state
}.

Fixpoint nfa_accept (x: A) w :=
match w with
| [] => nfa_fin A x
| a :: w => existsb y, (nfa_step A x a y) && nfa_accept y w
end.

Definition nfa_lang := [ pred w | nfa_accept (nfa_s A) w ].

```

The acceptance criterion given here avoids the matter of runs. In many cases, this will help us with proofs by induction on the accepted word. However, we will need runs in some of the proofs. Due to the fact that runs are not unique on NFAs, we will speak of **labeled paths** in this context. Labeled paths are paths through the graph implied by the automaton's step function. They can be valid runs for a word, but they could also be invalid by themselves in that they contain a sub-sequence of states that are not connected by edges. Therefore, we give a predicate that decides if a labeled path is a valid run for w on A . We then show that the acceptance criterion given above corresponds to the mathematical definition.

```

Fixpoint nfa_lpath x (xs : seq A) (w: word) {struct xs} :=
match xs, w with
| y :: xs', a :: w' => nfa_step A x a y && nfa_lpath y xs' w'
| [], [] => true
| -, _ => false
end.

Lemma nfa_lpath_accept x xs w:
  nfa_lpath x xs w -> nfa_fin A (last x xs) -> nfa_accept x w.

Lemma nfa_accept_lpath x w:
  nfa_accept x w -> exists xs, nfa_lpath x xs w /\ nfa_fin A (last x xs).

```

4.1.2 Deterministic Finite Automata

For functional δ , we speak of **deterministic** finite automata (DFA). In this case, we write δ as a function in our COQ development.

Listing 4.2: Deterministic Finite Automata

```

Record dfa : Type :=
{
  dfa_state :> finType;
  dfa_s : dfa_state;

```

```

    dfa_fin : pred dfa_state ;
    dfa_step : dfa_state -> char -> dfa_state
  }.
Fixpoint dfa_accept x w :=
match w with
| [] => dfa_fin A x
| a :: w => dfa_accept (dfa_step A x a) w
end.
Definition dfa_lang := [pred w | dfa_accept (dfa_s A) w].

```

Again, we avoid runs in our formalization of the acceptance criterion. In this case, however, we can give a function that computes the unique run of a word on A . This allows us to give an alternative acceptance criterion that is closer to the mathematical definition. We also prove that both criteria are equivalent.

```

Fixpoint dfa_run' (x: A) (w: word) : seq A :=
match w with
| [] => []
| a :: w => (dfa_step A x a) :: dfa_run' (dfa_step A x a) w
end.
Lemma dfa_run_accept x w: last x (dfa_run' x w) \in dfa_fin A = dfa_accept x w.

```

Equivalence between DFA and NFA

Deterministic and non-deterministic finite automata are equally expressive. One direction is trivial since every DFA can be seen as a NFA. We prove the other direction using the powerset construction. Given NFA A , we construct an equivalent DFA A_{det} in the following way:

$$\begin{aligned}
 Q_{det} &= \{q \mid q \subseteq Q\} \\
 s_{det} &= \{s\} \\
 F_{det} &= \{q \mid q \in Q_{det} \wedge q \cap F \neq \emptyset\} \\
 \delta_{det} &= \{(p, a, \bigcup_{p \in P} \{q \mid (p, a, q) \in \delta\}) \mid p, q \in Q_{det}, a \in \Sigma\}. \\
 A_{det} &= (Q_{det}, s_{det}, F_{det}, \delta_{det}).
 \end{aligned}$$

The formalization of A_{det} is straight-forward. We leave the set of states implicit.

```

Definition nfa_to_dfa :=
{[ dfa_s := set1 (nfa_s A);
  dfa_fin := [ pred X: {set A} | existsb x: A, (x \in X) && nfa_fin A x];
  dfa_step := [ fun X a => \bigcup (x | x \in X) finset (nfa_step A x a) ]
]}

```

Lemma 4.1.1. *For all powerset states X and for all states x with $x \in X$ we have that*

$$\mathcal{L}_x(A) \subseteq \mathcal{L}_X(A_{det}).$$

Proof. Let $w \in \mathcal{L}_x(A)$. We proof by induction on w that $w \in \mathcal{L}_X(A_{det})$.

- For $w = \varepsilon$ and $\varepsilon \in \mathcal{L}_x(A)$ we get $x \in F$ from $\varepsilon \in \mathcal{L}_x(A)$. From $x \in X$ we get $X \cap F \neq \emptyset$ and therefore $\varepsilon \in \mathcal{L}_X(A_{det})$.
- For $w = aw'$ and $aw' \in \mathcal{L}_x(A)$ we get y such that $w' \in \mathcal{L}_y(A)$ and $(x, a, y) \in \delta$. The latter gives us $y \in Y$ where Y is such that $(X, a, Y) \in \delta_{det}$. We apply the inductive hypothesis to y and Y and $w' \in \mathcal{L}_y(A)$, which gives us $w' \in \mathcal{L}_Y(A_{det})$. With $(X, a, Y) \in \delta_{det}$ we get $aw' \in \mathcal{L}_X(A_{det})$.

□

Lemma 4.1.2. *For all powerset states X and all words $w \in \mathcal{L}_X(A_{det})$ there exists a state x such that*

$$x \in X \wedge w \in \mathcal{L}_x(A).$$

Proof. We proof by induction on w that there exists x such that $x \in X \wedge w \in \mathcal{L}_x(A)$.

- For $w = \varepsilon$ and $\varepsilon \in \mathcal{L}_X(A_{det})$ we get $X \cap F \neq \emptyset$. Therefore, there exists x such that $x \in X$ and $x \in F$. Thus, we have $\varepsilon \in \mathcal{L}_x(A)$.
- For $w = aw'$ and $aw' \in \mathcal{L}_X(A_{det})$ we get Y such that $w' \in \mathcal{L}_Y(A_{det})$ and $(X, a, Y) \in \delta_{det}$. From the induction hypothesis we get y such that $y \in Y$ and $w' \in \mathcal{L}_y(A)$. From $y \in Y$ and $(X, a, Y) \in \delta_{det}$ we get x such that $x \in X$ and $(x, a, y) \in \delta$. Thus, $aw' \in \mathcal{L}_x(A)$.

□

Theorem 4.1.1. *The powerset automaton A_{det} accepts the same language as A , i.e.*

$$\mathcal{L}(A) = \mathcal{L}(A_{det}).$$

Proof. “ \subseteq ” This follows directly from lemma 4.1.1 with $x = s$ and $X = s_{det}$.

“ \supseteq ” From lemma 4.1.2 with $X = s_{det} = \{s\}$ we get $\mathcal{L}_{s_{det}}(A_{det}) \subseteq \mathcal{L}_s(A)$, which proves the claim. □

The formalization of this proof is straight-forward and follows the plan laid out above. The corresponding Lemmas are:

Lemma `nfa_to_dfa_complete` (x : A) w (X : nfa_to_dfa):

`x \in X -> nfa_accept A x w -> dfa_accept nfa_to_dfa X w.`

Lemma `nfa_to_dfa_sound` (X : nfa_to_dfa) w :

`dfa_accept nfa_to_dfa X w -> existsb x, (x \in X) && nfa_accept A x w.`

Lemma `nfa_to_dfa_correct` w : `nfa_lang A w = dfa_lang nfa_to_dfa w.`

4.2 Connected Components

Finite automata can have isolated subsets of states that are not reachable from the starting state. These states can not contribute to the language of the automaton. It will later be useful to have automata that only contain reachable states. Therefore, we define a procedure to extract the connected component from a given automaton.

Definition 4.2.1. Let $A = (\Sigma, Q, s, F, \delta)$ be a DFA. We define `reachable1` such that for all x and y , $y \in \text{reachable}(x) \iff \exists a, (x, a, y) \in \delta$. We write `reachable` to denote the transitive closure of `reachable1(s)`.

$$\begin{aligned} Q_c &= Q \cap \text{reachable} \\ s_c &= s \\ F_c &= F \cap \text{reachable} \\ \delta_c &= \{(x, a, y) \mid (x, a, y) \in \delta \wedge x, y \in Q_c\} \\ A_c &= (Q_c, s_c, F_c, \delta_c). \end{aligned}$$

We make use of SSREFLECT's `connect` predicate to extract a sequence of all states reachable from s . From this, we construct a finite type and use that as the new set of states. These new states carry a proof of reachability (`ssvalP`). We also have to construct a new transition function that ensures transitions always end in reachable states.

Definition `reachable1` := [fun x y => **existsb** a, dfa_step A1 x a == y].

Definition `reachable` := enum (connect reachable1 (dfa_s A1)).

Lemma `reachable0` : dfa_s A1 \in reachable.

Lemma `reachable_step` x a : x \in reachable -> dfa_step A1 x a \in reachable.

Definition `dfa_connected` :=

```
{|
  dfa_s := {|ssvalP := reachable0|};
  dfa_fin := [fun x => match x with {|ssval := x|} => dfa_fin A1 x end];
  dfa_step := [fun x a => match x with
    | {|ssvalP := Hx|} => {| ssvalP := (reachable_step _ a Hx) |}
  end]
|}.

```

Lemma 4.2.1. For every state $x \in \text{reachable}$ we have that

$$\mathcal{L}_x(A_c) = \mathcal{L}_x(A).$$

Proof. “ \subseteq ” Trivial. “ \supseteq ” We do an induction on w .

- For $w = \varepsilon$ we have $\varepsilon \in \mathcal{L}_x(A)$ and therefore $x \in F$. With $x \in \text{reachable}$ we get $x \in F_c$. Thus, $\varepsilon \in \mathcal{L}_x(A_c)$.

- For $w = aw'$ we have $y \in Q$ such that $(x, a, y) \in \delta$ and $w' \in \mathcal{L}_y(A)$. From $x \in \text{reachable}$ we get $y \in \text{reachable}$ by transitivity. Therefore, $(x, a, y) \in \delta_c$. The induction hypothesis gives us $w' \in \mathcal{L}_y(A_c)$. Thus, $aw' \in \mathcal{L}_x(A_c)$.

□

Theorem 4.2.1. *The language of the connected automaton A_c is identical to that of the original automaton A , i.e.*

$$\mathcal{L}(A) = \mathcal{L}(A_c).$$

Proof. By reflexivity, we have $s \in \text{reachable}$. We use lemma 4.2.1 with $x = s$ to prove the claim. □

The formalization of lemma 4.2.1 and theorem 4.2.1 is straight-forward.

Lemma `dfa_connected_correct' × (Hx: x \in reachable) :`
`dfa_accept dfa_connected {|ssvalP := Hx|} =1 dfa_accept A1 x.`

Lemma `dfa_connected_correct: dfa_lang dfa_connected =1 dfa_lang A1.`

To make use of the fact that A_c is fully connected, we will proof the characteristic property of A_c .

Definition 4.2.2. *A **representative** of a state x is a word w such that the unique run of w on A_c ends in x .*

Lemma 4.2.2. *We can give a representative for every state $x \in Q_c$.*

Sigma?

Proof. x carries a proof of reachability. This is equivalent to a path through the graph implied by `reachable1` that ends in x . We build the representative by extracting the edges taken in the path and building a word from those.

□

The formalization of theorem 4.2.2 consists of a more general version of the theorem, which facilitates the proof by induction over the path, and the theorem itself.

Lemma `dfa_connected_repr' (x y: dfa_connected):`
`connect reachable1_connected y x ->`
`exists w, last y (dfa_run' dfa_connected y w) = x.`

Lemma `dfa_connected_repr x :`
`exists w, last (dfa_s dfa_connected) (dfa_run dfa_connected w) = x.`

4.3 Emptiness

Given an automaton A , we can check if $\mathcal{L}(A) = \emptyset$. We simply obtain the connected automaton of A and check if there are any final states left.

Lemma 4.3.1. *The language of the connected automaton A_c is empty if and only if its set of final states F_c is empty, i.e.*

$$\mathcal{L}(A_c) = \emptyset \iff F_c = \emptyset.$$

Proof. “ \Leftarrow ” We have $\mathcal{L}(A_c) = \emptyset$ and have to show that for all $x \in Q_c$, $x \notin F_c$. Let $x \in Q_c$. By lemma 4.2.2 we get w such that the unique run of w on A_c ends in x . We use $\mathcal{L}(A_c) = \emptyset$ to get $w \notin \mathcal{L}(A_c)$, which implies that the run of w on A_c ends in a non-final state. By substituting the last state of the run by x we get $x \notin F_c$. “ \Rightarrow ” We have $F_c = \emptyset$ and have to show that for all words w , $w \notin \mathcal{L}(A_c)$. We use $F_c = \emptyset$ to show that the last state of the run of w on A_c is non-final. Thus, $w \notin \mathcal{L}(A_c)$. \square

Theorem 4.3.1. *We can decide emptiness of $\mathcal{L}(A)$ by computing the cardinality of A_c ’s set of final states, i.e.*

$$F_c = \emptyset \iff \mathcal{L}(A) = \emptyset.$$

Proof. We use theorem 4.2.1 to substitute $\mathcal{L}(A)$ by $\mathcal{L}(A_c)$. We then use lemma 4.3.1 to prove the claim. \square

The formalization of lemma 4.3.1 is split in two parts to facilitate its application.

Definition `dfa_lang_empty := #|dfa_fin dfa_connected| == 0.`

Lemma `dfa_lang_empty_complete: dfa_lang dfa_connected =1 pred0 -> dfa_lang_empty.`

Lemma `dfa_lang_empty_sound: dfa_lang_empty -> dfa_lang dfa_connected =1 pred0.`

Lemma `dfa_lang_empty_correct: dfa_lang_empty <-> dfa_lang A1 =1 pred0.`

4.4 Deciding Equivalence of Finite Automata

Definition 4.4.1. *We say that two automata are **equivalent** if and only if their languages are equal.*

Given finite automata A_1 and A_2 , we construct DFA A such that the language of A is the symmetric difference of the languages of A_1 and A_2 , i.e.,

$$\mathcal{L}(A) = \mathcal{L}(A_1) \ominus \mathcal{L}(A_2) = \mathcal{L}(A_1) \cap \neg \mathcal{L}(A_2) \cup \mathcal{L}(A_2) \cap \neg \mathcal{L}(A_1).$$

Theorem 4.4.1. *The equivalence of A_1 and A_2 is decidable, i.e.*

$$\mathcal{L}(A_1) = \mathcal{L}(A_2) \text{ if and only if } \mathcal{L}(A) \text{ is empty.}$$

Proof. The correctness of this procedure follows from the properties of the symmetric difference operator, i.e.

$$\mathcal{L}(A_1) \ominus \mathcal{L}(A_2) = \emptyset \Leftrightarrow \mathcal{L}(A_1) = \mathcal{L}(A_2).$$

The decidability of this procedure follows directly from theorem 4.3.1. \square

Listing 4.3: Formalization of theorem 4.4.1

Definition `dfa_sym_diff` :=
`dfa_disj (dfa_conj A1 (dfa_compl A2)) (dfa_conj A2 (dfa_compl A1)).`
Lemma `dfa_sym_diff_correct`:
`dfa_lang dfa_sym_diff =1 pred0 <=> dfa_lang A1 =1 dfa_lang A2.`

4.5 Regular Expressions and Finite Automata

We prove that there is a finite automaton for every extended regular expression and vice versa. In fact, we can give a standard regular expression for every finite automaton. With this, we will prove that extended regular expressions are equivalent to standard regular expressions, thereby proving closure under intersection and negation.

4.5.1 Regular Expressions to Finite Automata

We prove that there exists an equivalent automaton for every extended regular expressions. The structure of this proof is given by the inductive nature of regular expressions. For every constructor, we provide an equivalent automaton.

Include
all
proofs?

4.5.2 Deciding Equivalence of Regular Expressions

Based on our procedure to construct an equivalent automaton from a regular expression, we can decide equivalence of regular expressions. Given r_1 and r_2 , we construct equivalent DFA A_1 and A_2 as above.

4.5.3 Finite Automata to Regular Expressions

We prove that there is an equivalent standard regular expression for every finite automaton.

Since we are given an automaton it is not obvious how to partition our proof obligations into smaller parts. We use Kleene's original proof, the transitive closure method. This method recursively constructs a regular expression that is equivalent to the given automaton. Given a DFA A , we first assign some ordering to its states. We then define $R_{i,j}^k$ such that $\mathcal{L}(R_{i,j}^k)$ is the set of all words that have a run on A starting in state i that ends in state j without ever leaving a state smaller than k . The base case $R_{i,j}^0$ is the set of all singleton words that are edges between state i and j , and ε if $i = j$. Given $R_{i,j}^k$ we can easily define $R_{i,j}^{k+1}$ based on the observation that only one new state has to be considered:

Insert complete formal definition

$$R_{i,j}^{k+1} = R_{i,k}^k \cdot (R_{k,k}^k)^* \cdot R_{k,j}^k + R_{i,j}^k.$$

We make use of SSREFLECT's ordinals to get an ordering on states. We chose to employ ordinals for i and j , but not for k . This simplifies the inductive definitions on k . It does, however, lead to explicit conversions when k is used in place of i or j . In fact, i and j are states in our COQ implementation. We only rely on ordinals for comparison to k .

Add implementation of R

Furthermore, we define $L_{i,j}^k \subseteq \mathcal{L}(A)$ in terms of runs on the automaton. The relation of $L_{i,j}^k$ to $\mathcal{L}(A)$ can be proven very easily. We will also prove it equivalent to $R_{i,j}^k$. This allows us to connect $R_{i,j}^k$ to $\mathcal{L}(A)$.

Definition `allbutlast` $p : \text{pred } (\text{seq } X) :=$
`fun xs => all p (belast xs).`

Definition `L` $:=$
`[fun k: nat =>`
`[fun x y: A =>`
`[pred w |`
`(last x (dfa_run' A x w) == y)`
`&& allbutlast (<.k) (dfa_run' A x w)`
`]`
`]`
`].`

Theorem 4.5.1. *We can express $\mathcal{L}(A)$ in terms of L . L is equivalent to R .*

$$\mathcal{L}(A) = \bigcup_{f \in F} L_{s,f}^{|Q|} = \mathcal{L}\left(\sum_{f \in F} R_{s,f}^{|Q|}\right).$$

Proof. By definition, every $w \in \mathcal{L}(A)$ has a run that ends in some $f \in F$. Then, by definition, $w \in L_{s,f}^{|Q|}$.

It remains to show that $\mathcal{L}(R_{i,j}^k) = L_{i,j}^k$. This claim can be proven by induction over k . We begin with the inclusion of $\mathcal{L}(R_{i,j}^k)$ in $L_{i,j}^k$. For $k = 0$, we do a case distinction on $i == j$ and unfold R . The resulting three cases ($i == j \wedge w = \varepsilon$, $i == j \wedge |w| = 1$, $i <> j \wedge |w| = 1$) are easily closed.

The inductive step has two cases: A triple concatenation and a simple recursion. The second case is solved by the inductive hypothesis. In the first case, we split up the concatenation such that

$$w = w_1 \cdot w_2 \cdot w_3 \wedge w_1 \in \mathcal{L}(R_{i,k}^k) \wedge w_2 \in \mathcal{L}((R_{k,k}^k)^*) \wedge w_3 \in \mathcal{L}(R_{k,j}^k).$$

The induction hypothesis is applied to w_1 and w_3 to get $w_1 \in L_{i,k}^k$ and $w_3 \in L_{k,j}^k$. We use a lemma by Coquand and Siles that splits w_2 into a sequence of words from $\mathcal{L}(R_{k,k}^k)$ to which we can apply the induction hypothesis. Two concatenation lemmas for L are used to merge the sequence of words proven to be in $L_{k,k}^k$, w_1 and w_3 . This shows $\mathcal{L}(R_{i,j}^k) \subseteq L_{i,j}^k$.

Next, we show the inclusion of $L_{i,j}^k$ in $\mathcal{L}(R_{i,j}^k)$, again by induction over k . The base case is solved by case distinction on $i == j$. The inductive step requires a **splitting lemma** for L which shows that every non-empty word in $L_{i,j}^{k+1}$ is either in $L_{i,j}^k$ or has a non-empty prefix in $L_{i,k}^k$ and a corresponding suffix in $L_{k,j}^{k+1}$. In the first case, we can apply the induction hypothesis. In the second case, we use size induction on the word, apply the original induction hypothesis to the prefix and the size induction hypothesis to the suffix. We use two concatenation lemmas for R to merge the sub-expression. This finishes the proof. \square

Formalizing theorem 4.5.1 requires infrastructure to deal with *allbutlast*. Once this is in place, we can formalize the concatenation lemmas for R and L . These are required later to connect sub-results.

Lemma R_catL k i j $w1$ $w2$:

$$\begin{aligned} w1 \setminus in R^k i \ (k_ord \ k) &\rightarrow \\ w2 \setminus in R^{k+1} (k_ord \ k) \ j &\rightarrow \\ w1++w2 \setminus in R^{k+1} i \ j. \end{aligned}$$

Lemma L_catL k i j $w1$ $w2$:

$$\begin{aligned} w1 \setminus in L^k i \ (enum_val \ (k_ord \ k)) &\rightarrow \\ w2 \setminus in L^{k+1} (enum_val \ (k_ord \ k)) \ j &\rightarrow \\ w1++w2 \setminus in L^{k+1} i \ j. \end{aligned}$$

Lemma L_catL k i j $w1$ $w2$:

$$\begin{aligned} w1 \setminus in L^k i \ (enum_val \ (k_ord \ k)) &\rightarrow \\ w2 \setminus in L^{k+1} (enum_val \ (k_ord \ k)) \ j &\rightarrow \\ w1++w2 \setminus in L^{k+1} i \ j. \end{aligned}$$

We also need the splitting lemma mentioned earlier. This is quite intricate. We could split right after the first character and thereby simplify the lemma. However, the current form has the advantage of requiring simple concatenation lemmas.

Lemma `L_split k' i j a w:`
`let k := k_ord k' in`
`(a::w) \in L^k'.+1 i j ->`
`(a::w) \in L^k' i j /\`
`exists w1, exists w2,`
`a::w = w1 ++ w2 /\`
`w1 != [] /\`
`w1 \in L^k' i (enum_val k) /\`
`w2 \in L^k'.+1 (enum_val k) j.`

These lemmas suffice to show the claim of theorem 4.5.1.

Lemma `R.L_star k vv:`
`(forall (i j : 'L_#|A|) (w : word char),`
`w \in R^k i j -> w \in L^k (enum_val i) (enum_val j)) ->`
`all [predD mem_reg (R^k (k_ord k) (k_ord k)) &`
`eps (symbol:=char)] vv ->`
`flatten vv \in L^k.+1 (enum_val (k_ord k)) (enum_val (k_ord k)).`

Lemma `R.L k i j w:` `w \in R^k i j -> w \in L^k (enum_val i) (enum_val j).`

Lemma `L.R_1 k i j w:`
`(forall (i j : 'L_#|A|) (w : automata.word char),`
`w \in L^k (enum_val i) (enum_val j) -> w \in R^k i j) ->`
`w \in L^k.+1 (enum_val i) (enum_val j) -> w \in R^k.+1 i j.`

Lemma `L.R k i j w:` `w \in L^k (enum_val i) (enum_val j) -> w \in R^k i j.`

Fix this
mess

Chapter 5

Myhill-Nerode

The last characterization we consider is given by the Myhill-Nerode theorem.

5.1 Definition

The following definitions (taken from [7]) will lead us to the statement of the Myhill-Nerode theorem. We assume \equiv to be an equivalence relation on Σ^* , and L to be a language over Σ .

- (i) \equiv is **right congruent** if and only if for all $x, y \in \Sigma^*$ and $a \in \Sigma$,

$$x \equiv y \Rightarrow x \cdot a \equiv y \cdot a.$$

- (ii) \equiv **refines** L if and only if for all $x, y \in \Sigma^*$,

$$x \equiv y \Rightarrow (x \in L \Leftrightarrow y \in L).$$

- (iii) \equiv is of **finite index** if and only if it has finitely many equivalence classes, i.e.

$$\{[x] \mid x \in \Sigma^*\} \text{ is finite}$$

Definition 5.1.1. A relation is Myhill-Nerode if and only if it satisfies properties (i), (ii) and (iii).

Fix everything below this line

Given a language L , the Myhill-Nerode relation \approx_L is defined such that

$$\forall u, v \in \Sigma^*. u \approx_L v \iff \forall w \in \Sigma^*. u \cdot w \in L \Leftrightarrow v \cdot w \in L.$$

Listing 5.1: Myhill-Nerode relation

Definition $MN \ w1 \ w2 := \text{forall } w3, w1++w3 \ \backslash \text{in } L == (w2++w3 \ \backslash \text{in } L).$

Theorem 5.1.1. Myhill-Nerode Theorem. A language L is regular if and only if \approx_L is of finite index.

5.2 Finite Partitionings and Equivalence Classes

CoQ does not have quotient types. We pair up functions and proofs for certain properties of those functions to emulate quotient types.

A finite partitioning is a function from Σ^* to some finite type F . We use this concept to model equivalent classes in CoQ. A finite partitioning of the Myhill-Nerode relation is a finite partitioning f that also respects the Myhill-Nerode relation, i.e.,

$$\forall u, v \in \Sigma^*. f(u) = f(v) \Leftrightarrow u \approx_L v.$$

Listing 5.2: Finite partitioning of the Myhill-Nerode relation

Definition `MN_rel (f: Fin_eq_cls) := forall w1 w2, f w1 == f w2 <-> MN w1 w2.`

Theorem 5.2.1. *\approx_L is of finite index if and only if there exists a finite partitioning of the Myhill-Nerode relation.*

Proof. If \approx_L is of finite index, we use the set equivalence classes as a finite type and construct f such that

$$\forall w. f(w) = [w]_{\approx}.$$

f is a finite partitioning of the Myhill-Nerode relation by definition.

Conversely, if we have a finite partitioning of the Myhill-Nerode relation, we can easily see that \approx_L must be of finite index since f 's values directly correspond to equivalence classes. The image of f is finite. Therefore, \approx_L is of finite index. \square

A more general concept is that of a refining finite partitioning of the Myhill-Nerode relation:

$$\forall u, v \in \Sigma^*. f(u) = f(v) \Rightarrow u \approx_L v.$$

Listing 5.3: Refining finite partitioning of the Myhill-Nerode relation

Definition `MN_ref (f: Fin_eq_cls) := forall w1 w2, f w1 == f w2 -> MN w1 w2.`

We require all partitionings to be surjective. Therefore, every equivalence class x has at least one class representative which we denote $cr(x)$. Mathematically, this is not a restriction since there are no empty equivalence classes. In our constructive setting we would have to give a procedure that builds a minimal finite type F' from F and a corresponding function f' from Σ^* to F' such that f' is surjective and extensionally equal to f .

5.3 Minimizing Equivalence Classes

We will prove that refining finite partitionings can be converted into finite partitionings. For this purpose, we employ the table-filling algorithm to find indistinguishable states under the Myhill-Nerode relation ([5]). However, we do not rely on an automaton. In fact, we use the finite type F , i.e., the equivalence classes, instead of states.

Given a refining finite partitioning f , we construct a fixed-point algorithm. The algorithm initially outputs the set of equivalence classes that are distinguishable by the inclusion of their class representative in L . We denote this initial set $dist_0$.

$$dist_0 := \{(x, y) \in F \times F \mid cr(x) \in L \Leftrightarrow cr(y) \notin L\}.$$

To find more distinguishable equivalence classes, we have to identify equivalence classes that lead to distinguishable equivalence classes.

Definition 5.3.1. We say that a pair of equivalence classes (x, y) *transitions* to (x', y') with a if and only if

$$f(cr(x) \cdot a) = x' \wedge f(cr(y) \cdot a) = y'.$$

We denote (x', y') by $ext_a(x, y)$.

The fixed-point algorithm tries to extend the set of distinguishable equivalence classes by looking for a so-far undistinguishable pair of equivalence classes that transitions to a pair of distinguishable equivalence classes.

Definition 5.3.2.

$$unnamed(dist) := dist_0 \cup dist \cup \{(x, y) \mid \exists a. ext_a(x, y) \in dist\}$$

Lemma 5.3.1. *unnamed is monotone and has a fixed-point.*

Proof. Monotonicity follows directly from the monotonicity of \cup . The number of sets in $F \times F$ is finite. Therefore, *unnamed* has a fixed point. □

Let **distinct** be the fixed point of *unnamed*. Let **equiv** be the complement of **distinct**.

Theorem 5.3.1. *f_{min} is a finite partitioning of the Myhill-Nerode relation on L .*

Finish
construc-
tion

Add for-
maliza-
tion

5.4 Finite Automata and Myhill-Nerode

We prove theorem 5.1.1 by proving it equivalent to the existence of an automaton that accepts L .

5.4.1 Finite Automata to Myhill-Nerode

Given DFA A , for all words w we define $f(w)$ to be the last state of the run of w on A .

Lemma 5.4.1. *f is a refining finite partitioning of the Myhill-Nerode relation on $\mathcal{L}(A)$.*

Proof. The set of states of A is finite. For all u, v and w we have that if $f(u) = f(v) = x$, i.e., the runs of u and v on A end in the exact same state x . From this, we get that for all w , runs of $u \cdot w$ and $v \cdot w$ on A also end in the same state. Therefore, $u \cdot w \in \mathcal{L}(A)$ if and only if $v \cdot w \in \mathcal{L}(A)$. \square

Theorem 5.4.1. *If L is accepted by DFA A , then there exists a finite partitioning of the Myhill-Nerode relation on L .*

Proof. From lemma 5.4.1 we get a refining finite partitioning f of the Myhill-Nerode relation on $\mathcal{L}(A)$. Since L is accepted by A , $L = \mathcal{L}(A)$. Therefore, f is a refining finite partitioning of the Myhill-Nerode relation on L . By theorem 5.3.1 there also exists a finite partition of the Myhill-Nerode relation on L . \square

5.4.2 Myhill-Nerode to Finite Automata

Chapter 6

Conclusion

Bibliography

- [1] Janusz A. Brzozowski. Derivatives of regular expressions. *J. ACM*, 11(4):481–494, 1964.
- [2] Thierry Coquand and Vincent Siles. A decision procedure for regular expression equivalence in type theory. In *CPP*, pages 119–134, 2011.
- [3] Ding-Shu Du and Ker-I Ko. *Problem Solving in Automata, Languages, and Complexity*. 2001.
- [4] Georges Gonthier, Assia Mahboubi, and Enrico Tassi. A Small Scale Reflection Extension for the Coq system. Rapport de recherche RR-6455, INRIA, 2008.
- [5] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to automata theory, languages, and computation - international edition (2. ed)*. Addison-Wesley, 2003.
- [6] S. C. Kleene. Representation of events in nerve nets and finite automata. *Automata Studies*, 1965.
- [7] Dexter Kozen. *Automata and computability*. Undergraduate texts in computer science. Springer, 1997.
- [8] Peter Linz. *An introduction to formal languages and automata (4. ed.)*. Jones and Bartlett Publishers, 2006.
- [9] The Coq development team. *The Coq proof assistant reference manual*. LogiCal Project, 2004. Version 8.0.
- [10] Anil Nerode. Linear automaton transformations. *Proceedings of the American Mathematical Society*, 9(4):541–544, ao 1958.