

Certified Parsing of Regular Languages

Denis Firsov and Tarmo Uustalu

Institute of Cybernetics, Tallinn University of Technology
Akadeemia tee 21, 12618 Tallinn, Estonia
{denis,tarmo}@cs.ioc.ee

Abstract. We report on a certified parser generator for regular languages using the Agda programming language. Specifically, we programmed a transformation of regular expressions into a Boolean-matrix based representation of nondeterministic finite automata (NFAs). And we proved (in Agda) that a string matches a regular expression if and only if the NFA accepts it. The proof of the if-part is effectively a function turning acceptance of a string into a parse tree while the only-if part gives a function turning rejection into a proof of impossibility of a parse tree.

1 Introduction

Parsing is the process of structuring a linear representation (sentence, a computer program, etc.) in accordance with a given grammar. Parsers are procedures which perform parsing. They are being used extensively in a number of disciplines: in computer science (for compiler construction, database interfaces, artificial intelligence), in linguistics (for text analysis, corpora analysis, machine translation, textual analysis of biblical texts), in typesetting chemical formulae, in chromosome recognition, and so on [4]. It is therefore clear that having correct parsers is important for all these disciplines. Surprisingly, relatively little research had been done in field of certified parsing.

However, with the recent development of programming languages with dependent types it has become possible to encode useful invariants in the types and prove properties of a program while implementing it. In this paper we use the system of dependent types of the Agda language [8], which is based on Martin-Löf's type theory. One of the basic ideas behind Martin-Löf's type theory is the Curry-Howard interpretation of propositions as types. A proposition is proved by writing a program of the corresponding type.

Dependent types allow types to talk about values. A classical example of a dependent type is the type of lists of a given length: $\text{Vec } A \ n$. Here A is the type of the elements and n is the length of the list. Having such a definition of vector, we can define “safe” functions. Let us look at definition of a safe head function:

$$\begin{aligned} \text{head} &: \forall \{A \ n\} \rightarrow \text{Vec } A \ (\text{suc } n) \rightarrow A \\ \text{head } (x :: xs) &= x \end{aligned}$$

This definition says that the function `head` accepts only vectors with at least one element. So, it is now the responsibility of the programmer to provide non-empty vectors, otherwise compilation will fail.

In this paper, we have adopted the same general technique of expressing properties of data in their types to implement a library of matrix operations, program a transformation of regular expressions into a Boolean matrix based representation of nondeterministic finite automata (NFAs) and prove it correct. The correctness proof turns NFAs effectively into parsers: the proof of acceptance or rejection of a string gives us a parse tree (a witness of matching) or a proof of impossibility of one.

Our Agda development can be found online at <http://cs.ioc.ee/~denis/cert-reg>.

2 Regular Expressions

Before we start, let us take care of the alphabet (Σ) of our interest. Since many functions require Σ as a parameter and an NFA must share its alphabet with the regular expressions, we define Σ as a global parameter for each module of the parser-generator library.

```
module ModuleName ( $\Sigma$  : Set) ( $\_ \stackrel{?}{=} \_$  : Decidable ( $\_ \equiv \_$  {A =  $\Sigma$ }))
```

It states that the module is parametrized by an alphabet and also a decidable equality on the alphabet.

The datatype for regular expressions is called `RegExp` and is defined as

```
data RegExp : Set where
   $\epsilon$       : RegExp
  ' _      :  $\Sigma \rightarrow$  RegExp
  _  $\cup$  _   : RegExp  $\rightarrow$  RegExp  $\rightarrow$  RegExp
  _  $\cdot$  _   : RegExp  $\rightarrow$  RegExp  $\rightarrow$  RegExp
  _  $^+$     : RegExp  $\rightarrow$  RegExp
```

The base cases are regular expressions for the empty string (ϵ) and for single-character strings (`' _`). The step cases are given by the regular operations: `_ \cup _` for union, `_ \cdot _` for concatenation and `_ $^+$` for iteration at least once. Note that instead of the Kleene star (`_ *`) we use plus (`_ $^+$`). This is more convenient for us and does not restrict generality, as star is expressible as choice between the empty string and plus.

Now, we need to specify when a string (an element of type `List Σ`) is in the language of the regular expression, which is also called matching. This is done by introducing a *parsing (or matching) relation* (denoted by `_ \blacktriangleright _`) between strings and regular expressions.

```
String : Set
String = List  $\Sigma$ 

data _  $\blacktriangleright$  _ : String  $\rightarrow$  RegExp  $\rightarrow$  Set where
  empty  : []  $\blacktriangleright$   $\epsilon$ 
  symb   : {x :  $\Sigma$ }  $\rightarrow$  [x]  $\blacktriangleright$  ' x
  unionl : {u : String} {e1 e2 : RegExp}  $\rightarrow$  u  $\blacktriangleright$  e1  $\rightarrow$  u  $\blacktriangleright$  e1  $\cup$  e2
  unionr : {u : String} {e1 e2 : RegExp}  $\rightarrow$  u  $\blacktriangleright$  e2  $\rightarrow$  u  $\blacktriangleright$  e1  $\cup$  e2
```

$$\begin{aligned}
\text{con} & : \{u v : \text{String}\} \{e_1 e_2 : \text{RegExp}\} \rightarrow u \blacktriangleright e_1 \rightarrow v \blacktriangleright e_2 \\
& \quad \rightarrow u \# v \blacktriangleright e_1 \cdot e_2 \\
\text{plus1} & : \{u : \text{String}\} \{e : \text{RegExp}\} \rightarrow u \blacktriangleright e \rightarrow u \blacktriangleright e^+ \\
\text{plus2} & : \{u v : \text{String}\} \{e : \text{RegExp}\} \\
& \quad \rightarrow u \blacktriangleright e \rightarrow v \blacktriangleright e^+ \rightarrow u \# v \blacktriangleright e^+
\end{aligned}$$

(Arguments enclosed in curly braces are implicit. The type checker will try to figure out the argument value for you. If the type checker cannot infer an implicit argument, then it must be provided explicitly, e.g., `symb {x}`.)

Let us now examine the constructors of the relation `__▶__`:

1. The constructor `empt` states that the empty string is in the language of the regular expression ϵ .
2. The constructor `symb` states that the string consisting of a single character `x` is in the language of the regular expression `'x'`.
3. The constructor `unionl` (`unionr`) states that if a string `u` is in the language defined by `e1` (`e2`), then `u` is also in the language of `e1 ∪ e2` (`e2 (e1)`).
4. The constructor `con` states that if a string `u` is in the language `e1` and a string `v` is in the language of `e2` then the concatenation of both strings is in the language of `e1 · e2`.
5. The constructor `plus1` states that if a string `u` is in the language of `e`, then it is also in the language of `e+`.
6. The constructor `plus2` states that if a string `u` is in the language of `e` and a string `v` is in the language of `e+`, then concatenation of both strings is in the language of `e+`.

Note that a proof that a string is in the parsing relation with a regular expression is a parse tree. Note that we do not introduce any notion of “raw” parse trees, a parse tree is always a parse tree of a specific string.

3 A Matrix Library

The transition relation of a nondeterministic finite automaton (NFA) can be viewed as a labeled directed graph. So it can be expressed as a family of incidence matrices (one matrix per label). In addition to the nice algebraic properties that this approach highlights, it allows us to compose automata (expressed with matrices) in various ways by using block operations.

We therefore formalize matrices and some important matrix operations and their properties.

3.1 Matrices and Matrix Operations

What sort of elements can a matrix contain? Our approach abstracts from the type of elements. But in order for matrix addition and multiplication to be well-defined and satisfy the standard properties, it must form a commutative semiring. In Agda we introduce a parametrized module

module Data.Matrix (sr : CommSemiRing)

This declaration says that the module `Data.Matrix` is parametrized by a commutative semiring `sr`. A semiring will be a record containing the carrier type `R`, the operations, and the proofs of the laws of the semiring.

Next we define a representation for matrices. A matrix is a vector of vectors, therefore the matrix type can be defined as

$$\begin{aligned} _ \times _ &: \mathbb{N} \rightarrow \mathbb{N} \rightarrow \text{Set} \\ k \times l &= \text{Vec} (\text{Vec } R \ l) \ k \end{aligned}$$

where `k` denotes the number of rows and `l` the number of columns in a matrix. Let us implement some of the most important operations on matrices:

Null Matrix. The zero or null matrix is a matrix with all entries zero.

$$\begin{aligned} \text{null} &: \{k \ l : \mathbb{N}\} \rightarrow k \times l \\ \text{null} &= \text{replicate} (\text{replicate } \text{zero}) \end{aligned}$$

Here `zero` is the additive identity element of the semiring. Note that the arguments `k` and `l` are implicit (enclosed in curly braces), so in the most cases we can omit them and the type checker will try to infer their values automatically.

Identity Matrix. The identity or unit matrix of size `k` is the `k × k` square matrix with ones on the main diagonal and zeros elsewhere.

$$\begin{aligned} \text{id} &: \{k : \mathbb{N}\} \rightarrow k \times k \\ \text{id } \{0\} &= [] \\ \text{id } \{\text{succ } k\} &= (\text{one} :: \text{replicate } \text{zero}) :: \text{zipWith } _ + _ \ \text{null } (\text{id } \{k\}) \end{aligned}$$

Here `one` is the multiplicative identity element.

Addition. Matrix addition is the operation of adding two matrices by adding corresponding entries together.

$$\begin{aligned} _ \oplus _ &: \{k \ l : \mathbb{N}\} \rightarrow k \times l \rightarrow k \times l \rightarrow k \times l \\ _ \oplus _ [] [] &= [] \\ _ \oplus _ (\text{rowA} :: A') (\text{rowB} :: B') &= \text{zipWith } _ + _ \ \text{rowA } \text{rowB} :: A' \oplus B' \end{aligned}$$

Note, that signature of addition requires the dimensions of the two input matrices to be equal.

Transposition. The transpose of a matrix `A` is another matrix $\langle A \rangle$ (or A^T in mathematical notation) created by writing the rows of `A` as the columns of $\langle A \rangle$.

$$\begin{aligned} \langle _ \rangle &: \{k \ l : \mathbb{N}\} \rightarrow k \times l \rightarrow l \times k \\ \langle [] \rangle &= \text{replicate } [] \\ \langle \text{rowA} :: A' \rangle &= \text{zipWith } _ :: _ \ \text{rowA } \langle A' \rangle \end{aligned}$$

Multiplication. Matrix multiplication is a binary operation that takes a pair of matrices and produces another matrix. If `A` is an `k × l` matrix and `B` is an `l × m` matrix, the result `A ⊗ B` of their multiplication is an `k × m` matrix

defined only if the number l of columns of the first matrix A is equal to the number of rows of the second matrix B .

```

_⊗_ : {k l m : ℕ} → k × l → l × m → k × m
_⊗_ [] = []
_⊗_ (rowA :: A') B = multRow :: A' ★ B
where
  multRow = map
    (λ colB → (foldr (_+_ ) zero (zipWith (_*_ ) rowA colB))) ⟨ B ⟩

```

The result of matrix multiplication is a matrix whose elements are found by multiplying the elements within a row from the first matrix by the associated elements within a column from the second matrix and summing the products.

In our library we have proved a number of basic properties of matrix transposition, addition, multiplication.

3.2 Block Operations

How to create matrices from smaller matrices systematically? This section describes an approach to block operations on matrices that has been advocated by Macedo and Oliveira [6]. It is centered around four operations corresponding to the injections and projections of the biproducts on the category of natural numbers and matrices.

- | | |
|--|--|
| <ol style="list-style-type: none"> 1. $\iota_1 : \{k \mid l : \mathbb{N}\} \rightarrow (k + l) \times k$
 $\iota_1 = \text{id} \# \text{null}$ 2. $\iota_2 : \{k \mid l : \mathbb{N}\} \rightarrow (k + l) \times l$
 $\iota_2 = \text{null} \# \text{id}$ | <ol style="list-style-type: none"> 3. $\pi_1 : \{k \mid l : \mathbb{N}\} \rightarrow k \times (k + l)$
 $\pi_1 = \text{zipWith } _ \# _ \text{id } \text{null}$ 4. $\pi_2 : \{k \mid l : \mathbb{N}\} \rightarrow l \times (k + l)$
 $\pi_2 = \text{zipWith } _ \# _ \text{null } \text{id}$ |
|--|--|

Example 1. Let us show some instances of these operations:

$$\begin{array}{cc}
 \iota_1\{3\}\{1\} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} & \iota_2\{2\}\{2\} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \\
 \pi_1\{3\}\{1\} = \begin{bmatrix} 1 & 0 & 0 & | & 0 \\ 0 & 1 & 0 & | & 0 \\ 0 & 0 & 1 & | & 0 \end{bmatrix} & \pi_2\{2\}\{2\} = \begin{bmatrix} 1 & 0 & | & 0 & 0 \\ 0 & 1 & | & 0 & 0 \end{bmatrix}
 \end{array}$$

The main block operations are now defined as follows:

Concatenation (copairing) is the operation of placing two matrices next to each other.

$$\begin{array}{l}
 [_ | _] : \{k \mid l \mid m : \mathbb{N}\} \rightarrow k \times l \rightarrow k \times m \rightarrow k \times (l + m) \\
 [A \mid B] = A \otimes \pi_1 \oplus B \otimes \pi_2
 \end{array}$$

Stacking (pairing) is the operation of placing two matrices on top of each other.

$$\begin{aligned} [_ / _] : \{k \mid m : \mathbb{N}\} &\rightarrow k \times m \rightarrow l \times m \rightarrow (k + l) \times m \\ [A / B] &= \iota_1 \otimes A \oplus \iota_2 \otimes B \end{aligned}$$

Example 2. Let $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ and $B = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$. Then

$$[A \mid B] = \begin{bmatrix} 1 & 2 & 5 & 6 \\ 3 & 4 & 7 & 8 \end{bmatrix} \text{ and } [A / B] = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{bmatrix}$$

3.3 Properties of Block Operations

The main advantage of working with injections and projections is that they structure proofs about block operations. Next, we list some properties of block operations.

- Multiplying the concatenation of A and B by ι_1 (ι_2) yields A (B).

$$m_1 m_2\text{-con-}\iota_1 : \{k \mid m : \mathbb{N}\} (A : k \times l) (B : k \times m)$$

$$\rightarrow [A \mid B] \otimes (\iota_1 \{l\} \{m\}) \equiv A$$

$$m_1 m_2\text{-con-}\iota_2 : \{k \mid m : \mathbb{N}\} (A : k \times l) (B : k \times m)$$

$$\rightarrow [A \mid B] \otimes (\iota_2 \{l\} \{m\}) \equiv B$$

- Multiplying π_1 (π_2) by the stacking of A and B gives A (B).

$$\pi_1\text{-}m_1 m_2\text{-stack} : \{k \mid m : \mathbb{N}\} (A : k \times l) (B : m \times l)$$

$$\rightarrow \pi_1 \{k\} \{m\} \otimes [A / B] \equiv A$$

$$\pi_2\text{-}m_1 m_2\text{-stack} : \{k \mid m : \mathbb{N}\} (A : k \times l) (B : m \times l)$$

$$\rightarrow \pi_2 \{k\} \{m\} \otimes [A / B] \equiv B$$

- The product of C and the concatenation of A and B is equal to the concatenation of $C \otimes A$ and $C \otimes B$.

$$\text{distrib-lft} : \{k \mid m \mid n : \mathbb{N}\} (A : k \times l) (B : k \times m) (C : n \times k)$$

$$\rightarrow C \otimes [A \mid B] \equiv [C \otimes A \mid C \otimes B]$$

- The product of a stacking reduces to a stacking of products.

$$\text{distrib-rgt} : \{k \mid m \mid n : \mathbb{N}\} (A : k \times l) (B : m \times l) (C : l \times n)$$

$$\rightarrow [A / B] \otimes C \equiv [A \otimes C / B \otimes C]$$

- Multiplying the concatenation of A and B by the stacking of C and D yields the sum of $A \otimes C$ and $B \otimes D$.

$$\text{con-}\otimes\text{-stack} : \{k \mid m \mid n : \mathbb{N}\} (A : k \times l) (B : k \times m)$$

$$\rightarrow (C : l \times n) \rightarrow (D : m \times n)$$

$$\rightarrow [A \mid B] \otimes [C / D] \equiv A \otimes C \oplus B \otimes D$$

4 NFAs and Parsing with NFAs

We are now in the position to implement a parser generator for regular languages. We parse strings with nondeterministic finite automata and represent them in terms of Boolean matrices.

From now on we therefore use matrices over the commutative semiring of Booleans, with false as zero, disjunction as addition, true as one, conjunction as multiplication.

4.1 Nondeterministic Finite Automata

A Σ -NFA can be defined as a record with four fields.

```

record NFA : Set where
  field  $\nabla$  :  $\mathbb{N}$ 
          $\delta$  :  $\Sigma \rightarrow \nabla \times \nabla$ 
          $l$  :  $1 \times \nabla$ 
          $F$  :  $\nabla \times 1$ 

```

- ∇ is the *size of the state space*. We do not name states. Instead we identify them with positions in rows and columns of matrices.
- δ specifies the *transition function*. In our implementation δ is a total function from letters of the alphabet to incidence matrices such that for any $x : \Sigma$ the function call δx will return an incidence matrix D of size $\nabla \times \nabla$ where $D_{ij} = 1$ iff q_j is a successor of q_i for character x .
- l specifies the *set of initial states*. The initial states can be represented by a $1 \times \nabla$ matrix (row vector) where the element l_{i1} is 1 iff q_i is an initial state.
- F specifies the *set of final states*. The final states can be represented by a $\nabla \times 1$ matrix (column vector) where the element F_{i1} is 1 iff q_i is a final state.

4.2 Running an NFA

Running an NFA on the string $x_0 \dots x_n$ from a set of states X represented by a row vector can be implemented as the series of multiplications $X \otimes \delta \text{ nfa } x_0 \otimes \dots \otimes \delta \text{ nfa } x_n$. This computes the row vector of all states reachable from the states X by following the transitions corresponding to the individual letters of the string $x_0 \dots x_n$.

```

run : (nfa : NFA)  $\rightarrow$  String  $\rightarrow$   $1 \times \nabla$  nfa  $\rightarrow$   $1 \times \nabla$  nfa
run nfa u X = foldl ( $\lambda A x \rightarrow A \otimes \delta \text{ nfa } x$ ) X u

```

If we take X to be l and multiply the matrix further with the column vector F of the final states, we get the 1×1 matrix $\text{id } \{1\}$ (i.e. $[[\text{true}]]$), if there is an overlap between the states reachable from some initial state and the final states, i.e., if the string is accepted, and null (i.e. $[[\text{false}]]$) otherwise.

```

runNFA : NFA  $\rightarrow$  String  $\rightarrow$   $1 \times 1$ 
runNFA nfa u = (run nfa u (l nfa))  $\otimes$  F nfa

```

4.3 Converting Regular Expressions to NFAs (Parsers)

We have introduced and defined the types `RegExp` and `NFA`. We now implement a conversion from `RegExp` to `NFA`, which we will use as a parser generator.

```

reg2nfa : RegExp → NFA
reg2nfa ε      = ε'
reg2nfa (' a)  = '' a
reg2nfa (e1 ∪ e2) = (reg2nfa e1) ∪' (reg2nfa e2)
reg2nfa (e+)    = (reg2nfa e)+'
reg2nfa (e1 · e2) = (reg2nfa e1) ·' (reg2nfa e2)
    
```

This function recurses over the regular expression and replaces every constructor with a corresponding operation on NFAs. We describe each case:

– $e = \varepsilon$:

```

ε' = record {
  ∇ = 1 ;
  δ = λ x → null;
  I = id;
  F = id}
    
```

Clearly, an NFA that accepts only the empty string can be given by one state 0 that is both initial and final.

– $e = ' a$:

In this case, the regular expression describes the single-character string a . So, the corresponding NFA should accept only this string.

```

'' a = record {
  ∇ = 2;
  δ = λ x → if a == x then [ [ null | id {1} ] /
                             [ null | null ] ]
                             else null;
  I = [ id {1} | null ];
  F = [ null / id {1} ]}
    
```

This NFA has two states 0 and 1 such that 0 is an initial and 1 a final state. The transition function compares the character x with the expected character a . If they coincide, then it returns the incidence matrix for the graph with a single edge from 0 to 1, otherwise, the null matrix for the empty graph.

– $e = e_1 \cup e_2$:

Recall that strings of both languages must be accepted. So we must run both NFAs.

```

nfa1 ∪' nfa2 = record {
  ∇ = ∇ nfa1 + ∇ nfa2;
  δ = λ x → [ [ δ nfa1 x | null ] /
              [ null | δ nfa2 x ] ];
  I = [ I nfa1 | I nfa2 ];
  F = [ F nfa1 / F nfa2 ]}
    
```


The resulting NFA is built of nfa_1 and nfa_2 as follows:

- (∇) The state space of the resulting NFA must contain states from nfa_1 and nfa_2 . So $\nabla = \nabla nfa_1 + \nabla nfa_2$.
- (δ) The transition function of the resulting NFA is composed of four blocks.
 1. The top left block is the incidence matrix of the first NFA.
 2. The top right block is null. So, no transitions from nfa_1 to nfa_2 .
 3. The bottom left block is null. In terms of incidence matrices this signals that there are no transitions from nfa_2 to nfa_1 .
 4. The bottom right block is the incidence matrix of the second NFA.
- (I) The set of initial states in the resulting NFA is the union of the sets of initial states of nfa_1 and nfa_2 .
- (F) The set of final states in the resulting NFA is the union of the sets of final states of nfa_1 and nfa_2 .

– $e = e^{+'}$:

$$\begin{aligned} nfa^{+'} = \mathbf{record} \{ \\ \nabla &= \nabla nfa; \\ \delta &= \lambda x \rightarrow (\text{id} \oplus F nfa \otimes I nfa) \otimes (\delta nfa x) \\ I &= I nfa; \\ F &= F nfa \} \end{aligned}$$

The difference between $nfa^{+'}$ and nfa is in δ only. Specifically, we add a new edge from each final state to each successor of an initial state. This is achieved by $F nfa \otimes I nfa \otimes \delta nfa x$, where $I nfa \otimes \delta nfa x$ stands for edges reachable from initial state by reading the token x . And $F nfa \otimes I nfa \otimes \delta nfa x$ puts an edge from each final state to each successor of an initial state.

– $e = e_1 \cdot e_2$:

$$\begin{aligned} nfa_1 \cdot nfa_2 = \mathbf{record} \{ \\ \nabla &= \nabla nfa_1 + \nabla nfa_2; \\ \delta &= \lambda x \rightarrow \left[\begin{array}{c|c} [\delta nfa_1 x \mid F nfa_1 \otimes I nfa_2 \otimes \delta nfa_2 x] & / \\ \text{null} & \delta nfa_2 x \end{array} \right]; \\ I &= [I nfa_1 \mid \text{null}]; \\ F &= [F nfa_1 \otimes I nfa_2 \otimes F nfa_2 \mid F nfa_2] \} \end{aligned}$$

The fields of nfa_1 and nfa_2 are combined in the following way:

- (∇) The state space of the resulting NFA consists of the disjoint union of state spaces of nfa_1 and nfa_2 , i.e. $\nabla = \nabla nfa_1 + \nabla nfa_2$.
- (δ) The transition function constructs incidence matrices from four blocks.
 1. The top left block contains the incidence matrix of nfa_1 . Hence, the transition relation between the states of nfa_1 is not changed.
 2. The top right block is $F nfa_1 \otimes I nfa_2 \otimes \delta nfa_2 x$. This expression constructs an incidence matrix with transitions from all final states of nfa_1 to all successors of initial states in nfa_2 . In other words, it says that upon reaching a final state of nfa_1 , it is time to transition to nfa_2 .
 3. The bottom left block is null. Hence, there are no transitions from nfa_2 back to nfa_1 .

4. The bottom right block consists of the transition function of nfa_2 . So, the transitions between the states of nfa_2 are unchanged.
- (I) The resulting NFA's initial states must contain only initial states of nfa_1 .
 - (F) Clearly, the final states of the resulting NFA must contain all final states of the second NFA. But what if the second NFA accepts the empty string? Then the resulting NFA must also accept the language of the first NFA, hence, contain all its final states. The desired behaviour is achieved by

$$F = [F\ nfa_1 \otimes I\ nfa_2 \otimes F\ nfa_2 \ / \ F\ nfa_2]$$

The top block of this column vector is equal to either the empty vector or $F\ nfa_1$ depending on the result of $I\ nfa_2 \otimes F\ nfa_2$. The latter multiplication is equal to $id\ \{1\}$, if nfa_2 accepts the empty string, and $null$, if it does not. The bottom block of F is always equal to $F\ nfa_2$. So, the desired behaviour is achieved.

4.4 Correctness

The correctness of an NFA with respect to a regular expression consists of completeness and soundness. Completeness guarantees that every string matching a regular expression will be accepted by the NFA. However, completeness alone is not enough, since an NFA accepting all strings is also complete. Soundness in turn guarantees that, if the NFA accepts, the string matches the regular expression. Similarly to the case of completeness, soundness is not sufficient alone, because an NFA rejecting every string is sound.

Completeness. Completeness states that, if a string is matched by the given regular expression, then the constructed NFA $reg2nfa\ e$ accepts it.

$$\begin{aligned} complete & : (e : RegExp) \rightarrow (u : String) \rightarrow u \blacktriangleright e \\ & \rightarrow runNFA (reg2nfa\ e)\ u \equiv id\ \{1\} \end{aligned}$$

We prove this theorem by induction on the proof of $u \blacktriangleright e$. Hence, all shapes of parse trees must be considered. We describe only the cases for union and concatenation, the others can be found in the Agda code.

Union In this case $e = e_1 \cup e_2$.

$$complete\ (e_1 \cup e_2)\ u\ parseTree = \dots$$

We start with pattern matching on $parseTree$.

$$\begin{aligned} complete\ (e_1 \cup e_2)\ u\ (unionl\ parseTree') & = \dots \\ complete\ (e_1 \cup e_2)\ u\ (unionr\ parseTree') & = \dots \end{aligned}$$

This yields two cases for the last rule used in the parse tree: `unionl` or `unionr`. Since both cases are proved in the same way, we describe only the first one.

The main idea is to show that a run of $nfa_1 \cup' nfa_2$ can be split into runs of nfa_1 and nfa_2 . It is proved by the lemma `union-split`:

$$\begin{aligned} \text{union-split} &: (\text{nfa}_1 \text{ nfa}_2 : \text{NFA}) (u : \text{String}) (X_1 : 1 \times _) (X_2 : 1 \times _) \\ &\rightarrow \text{run} (\text{nfa}_1 \cup' \text{nfa}_2) s [X_1 \mid X_2] \equiv [\text{run} \text{nfa}_1 s X_1 \mid \text{run} \text{nfa}_2 s X_2] \end{aligned}$$

Next, by using the previously described property

$$\begin{aligned} \text{con-}\otimes\text{-stack} &: \{k \mid m \ n : \mathbb{N}\} (A : k \times l) (B : k \times m) (C : l \times n) (D : m \times n) \\ &\rightarrow [A \mid B] \otimes [C / D] \equiv A \otimes C \oplus B \otimes D \end{aligned}$$

we complete the proof:

$$\frac{\frac{\text{id } \{1\} \oplus \text{id } \{1\} \equiv \text{id } \{1\}}{\text{Boolean arithm.}}}{\frac{(\text{run} \text{nfa}_1 s X_1 \otimes (\text{F} \text{nfa}_1)) \oplus (\text{run} \text{nfa}_2 s X_2 \otimes (\text{F} \text{nfa}_2)) \equiv \text{id } \{1\}}{\text{IHs}}}{\frac{[\text{run} \text{nfa}_1 s X_1 \mid \text{run} \text{nfa}_2 s X_2] \otimes [\text{F} \text{nfa}_1 / \text{F} \text{nfa}_2] \equiv \text{id } \{1\}}{\text{union-split}}}}{\text{run} (\text{nfa}_1 \cup' \text{nfa}_2) s [X_1 \mid X_2] \otimes [\text{F} \text{nfa}_1 / \text{F} \text{nfa}_2] \equiv \text{id } \{1\}}$$

Plus In this case $e = e'^+$.

$$\text{complete } (e'^+) u \text{ parseTree} = \dots$$

Pattern matching on `parseTree` yields two cases. We examine them in turn:

1. In the first case the last rule of the parse tree is `plus1`.

$$\text{complete } (e'^+) u (\text{plus1 } \text{parseTree}') = \dots$$

Recall that `plus1` is a constructor which states that, if u is in language of e' , then it is also in language of (e'^+) . Hence, the main lemma for this case can be stated as

$$\begin{aligned} \text{plus-weak} &: (\text{nfa} : \text{NFA}) (u : \text{String}) (X : 1 \times (\nabla \text{nfa})) \\ &\rightarrow \text{run} \text{nfa} u X \otimes \text{F} \text{nfa} \equiv \text{id } \{1\} \\ &\rightarrow \text{run} (\text{nfa}^+) u X \otimes \text{F} \text{nfa} \equiv \text{id } \{1\} \end{aligned}$$

It is proved by induction on the length of the string u .

2. In the second case the last rule of the parse tree is `plus2`.

$$\text{complete } (e'^+) . (u_1 \# u_2) (\text{plus2 } \{u_1\} \{u_2\} \text{tree}_1 \text{tree}_2) = \dots$$

Note that string u is now split into u_1 and u_2 such that u_1 is in the language of e' and u_2 is in the language of e'^+ . We must prove that $u_1 \# u_2$ are in the language of e'^+ . To do so, we first introduce some useful lemmas.

- We show that `run` on $u_1 \# u_2$ can be split into a `run` on u_1 and a `run` on u_2 .

$$\begin{aligned} \text{plus-split} &: (\text{nfa} : \text{NFA}) (u_1 \ u_2 : \text{String}) \\ &\rightarrow \text{run} (\text{nfa}^+) (u_1 \# u_2) (\text{l} \text{nfa}) \otimes \text{F} \text{nfa} \\ &\equiv \text{run} (\text{nfa}^+) u_2 (\text{run} (\text{nfa}^+) u_1 (\text{l} \text{nfa})) \otimes \text{F} \text{nfa} \end{aligned}$$

- We also show that, if the automaton nfa^+ accepts a string from the initial states, then it will also accept that string from any final state.

$$\begin{aligned} \text{plus-fin} &: (\text{nfa} : \text{NFA}) (u : \text{String}) (X : 1 \times \nabla \text{nfa}) \\ &\rightarrow \text{run} (\text{nfa}^+) u (\text{l} \text{nfa}) \otimes \text{F} \text{nfa} \equiv \text{id } \{1\} \\ &\rightarrow X \otimes \text{F} \text{nfa} \equiv \text{id } \{1\} \\ &\rightarrow \text{run} (\text{nfa}^+) u X \otimes \text{F} \text{nfa} \equiv \text{id } \{1\} \end{aligned}$$

Finally, the big picture of the proof looks like this:

$$\frac{\frac{\text{run } (nfa^{+'}) \ u_2 \ (l \ nfa) \ \otimes \ F \ nfa \ \equiv \ id \ \{1\}}{\text{run } (nfa^{+'}) \ u_1 \ (l \ nfa) \ \otimes \ F \ nfa \ \equiv \ id \ \{1\}} \text{ IH} \quad \frac{\text{run } nfa \ u_1 \ (l \ nfa) \ \otimes \ F \ nfa \ \equiv \ id \ \{1\}}{\text{run } (nfa^{+'}) \ u_1 \ (l \ nfa) \ \otimes \ F \ nfa \ \equiv \ id \ \{1\}} \text{ plus-weak}}{\frac{\text{run } (nfa^{+'}) \ u_2 \ (\text{run } (nfa^{+'}) \ u_1 \ (l \ nfa)) \ \otimes \ F \ nfa \ \equiv \ id \ \{1\}}{\text{run } (nfa^{+'}) \ (u_1 \ \# \ u_2) \ (l \ nfa) \ \otimes \ F \ nfa \ \equiv \ id \ \{1\}} \text{ plus-split}} \text{ plus-fin}$$

Soundness. Showing that our NFA generation is sound is more complicated, but also more interesting. Let us look at the signature of the soundness theorem:

$$\begin{aligned}
 \text{sound} &: (e : \text{RegExp}) \rightarrow (u : \text{String}) \\
 &\rightarrow \text{runNFA } (\text{reg2nfa } e) \ u \ \equiv \ id \ \{1\} \rightarrow u \ \blacktriangleright \ e
 \end{aligned}$$

It states that, if the NFA accepts a string, then it matches the regular expression. `sound` is a proposition, but it is also a type! Its proof is a function that delivers parse trees. We prove this theorem by induction on the argument `e : RegExp`.

As in case of completeness, our aim is to explain the high-level ideas of the proof. We skip most of the details and describe only two cases.

Single character. This case is interesting because it demonstrates the essence of all soundness cases. We are given an accepting run of the automaton. Using this fact we must construct a parse tree. However, most of the cases generated by pattern matching are discharged by showing that they contradict with the accepting run we have at our disposal.

$$\text{sound } (' \ a) \ u \ \text{run} = \dots$$

We pattern match on the string `u` and examine three different cases in turns:

1. *u is the empty string.*

$$\text{sound } (' \ a) \ u \ \text{run} = \dots$$

We must show that in this case it is impossible to give `run`. We do so by pattern matching on `run` and the rest is taken care of by Agda's type checker.

2. *u is a string of one symbol.*

$$\text{sound } (' \ a) \ (x :: []) \ \text{run} = \dots$$

This is the only situation when the automaton can accept `u`. Still, we must check if `x` is equal to `a`. We case analyse on the decidable equality of `a` and `x`:

$$\text{sound } (' \ a) \ (x :: []) \ \text{run} \ \mathbf{with} \ a \stackrel{?}{=} x$$

$$\text{sound } (' \ a) \ (x :: []) \ | \ \text{eq} = \dots$$

$$\text{sound } (' \ a) \ (x :: []) \ | \ \text{neq} = \dots$$

Then two cases must be discharged.

- (a) *x is equal to a.*

This is exactly the case when the automaton finishes in the accepting state. To close this case, we rewrite the context using `a ≡ x` and provide the constructor `symb {x}` as the required proof.

(b) x is not equal to a .

Then Agda computes that $\text{runNFA } (\text{reg2nfa } (' a)) [x]$ is equal to null, but this contradicts the assumptions. Hence, the case is discharged.

3. u is a string of two or more characters.

$\text{sound } (' a) (x_1 :: x_2 :: xs) \text{ run} = \dots$

Similarly to the first case, our goal is to show that $\text{runNFA } (\text{reg2nfa } (' a)) u$ will never accept a string consisting of two or more characters. This is done by observing the fact that, even if the automaton reaches the second state by reading the first character, then by reading the second character the automaton will lose all active states, since there are no transitions going out of the second state. Hence, $\text{runNFA } (\text{reg2nfa } (' a)) (x_1 :: x_2 :: xs)$ cannot return $\text{id } \{1\}$. Therefore, the case is discharged.

Concatenation. We are in the case

$\text{sound } (e_1 \cdot e_2) u \text{ run} = \dots$

Fortunately, there is only one possible constructor for this case in the parsing relation, namely con . It states that to prove $u \blacktriangleright e_1 \cdot e_2$ we must show that $u_1 \blacktriangleright e_1$ and $u_2 \blacktriangleright e_2$ for some splitting of u into u_1 and u_2 . Hence, we must be able to extract from run two shorter runs and use them to get $u_1 \blacktriangleright e_1$ and $u_2 \blacktriangleright e_2$ by induction hypothesis. To express this in Agda we use sigma-types, corresponding to existentials.

```
cons-split : (nfa1 nfa2 : NFA) (u : String)
  → run (nfa1 · nfa2) u [ l nfa1 | null ]
    ⊗ [ F nfa1 ⊗ l nfa2 ⊗ F nfa2 / F nfa2 ] ≡ id {1}
  → ∃ [u1 : String] ∃ [u2 : String] u ≡ u1 ++ u2
    ∧ run nfa1 u1 (l nfa1) ⊗ F nfa1 ≡ id {1}
    ∧ run nfa2 u2 (l nfa2) ⊗ F nfa2 ≡ id {1}
```

Note that we split the string u into u_1 and u_2 , but must also provide a proof that $u \equiv u_1 ++ u_2$. To prove it, we will need a variant of cons-split .

```
cons-split-state : (nfa1 nfa2 : NFA) (x : Σ) (u : String) (X : 1 × _)
  → run (nfa1 · nfa2) (x :: u) [ X | null ] ⊗ (null ++ F nfa2) ≡ id {1}
  → ∃ [u1 : String] ∃ [u2 : String] (x :: u) ≡ u1 ++ u2
    ∧ run nfa1 u1 X ⊗ F nfa1 ≡ id {1}
    ∧ run nfa2 u2 (l nfa2) ⊗ F nfa2 ≡ id {1}
```

The important differences between cons-split and cons-split-state are the following:

- In cons-split , the given run starts from $[l nfa_1 | \text{null}]$, where $l nfa_1$ is the set of initial states of nfa_1 , but in cons-split-state , we use a more general variant $[X | \text{null}]$, where X is a given parameter.
- In cons-split , the run of the automaton can terminate either in the final states of nfa_1 or in the final states of nfa_2 , but in cons-split-state , the set of final states is limited to those of nfa_2 .

- In contrast with `cons-split-state`, the string `xs` can be empty in `cons-split`.

The restrictions present in `cons-split-state` force it to address the specific and most complicated case when the given run starts in some states of `nfa1`, but terminates in `nfa2`. The lemma states that in this case we can break an accepting run of `nfa1 ·' nfa2` down into two smaller accepting runs.

We will not describe here how `cons-split-state` is proved. Instead we show how to reduce `cons-split` to `cons-split-state`. To do so, we perform multiple case analyses. First, we distinguish the empty string case from the `cons`-case.

1. If $u \equiv []$, then $u_1 \equiv []$ and $u_2 \equiv []$ and we must show that $\text{I nfa}_1 \otimes \text{F nfa}_1 \equiv \text{id } \{1\}$ and $\text{I nfa}_2 \otimes \text{F nfa}_2 \equiv \text{id } \{1\}$. Both proofs are easily derived from the given premise:

$$[\text{I nfa}_1 \mid \text{null}] \otimes [\text{F nfa}_1 \otimes (\text{I nfa}_2 \otimes \text{F nfa}_2) / \text{F nfa}_2] \equiv \text{id } \{1\}$$

2. If $u \equiv x :: xs$, then we perform an additional case analysis: whether the second automaton has final states among its initial states:
 - (a) $\text{I nfa}_2 \otimes \text{F nfa}_2 \equiv \text{null}$. In this case, the problem is clearly an instance of `cons-split-state`.
 - (b) $\text{I nfa}_2 \otimes \text{F nfa}_2 \equiv \text{id } \{1\}$. In this case, we perform a third level of case analysis: whether the first automaton `nfa1` accepts the whole string.
 - i. $\text{run nfa}_1 u (\text{I nfa}_1) \otimes \text{F nfa}_1 \equiv \text{id } \{1\}$.
If we take $u_1 \equiv u$ and $u_2 \equiv []$, then case is immediately discharged.
 - ii. $\text{run nfa}_1 u (\text{I nfa}_1) \otimes \text{F nfa}_1 \equiv \text{null}$.

In this case, we need an additional lemma

$$\begin{aligned} \text{consnd2} &: (\text{nfa}_1 \text{ nfa}_2 : \text{NFA}) (u : \text{String}) \\ &\rightarrow (\text{run nfa}_1 u (\text{I nfa}_1)) \otimes \text{F nfa}_1 \equiv \text{null} \\ &\rightarrow (\text{run } (\text{nfa}_1 \cdot' \text{nfa}_2) u [\text{I nfa}_1 \mid \text{null}]) \otimes [\text{F nfa}_1 / \text{F nfa}_2] \\ &\equiv (\text{run } (\text{nfa}_1 \cdot' \text{nfa}_2) u [\text{I nfa}_1 \mid \text{null}]) \otimes [\text{null} / \text{F nfa}_2] \end{aligned}$$

It states that, if the first automaton does not accept the whole string, then running the automaton `nfa1 ·' nfa2` with the final states of `nfa1` is equivalent to running the automaton `nfa1 ·' nfa2` without the final states of `nfa1`. So, this branch of `cons-split` is also reduced to `cons-split-state`.

`sound` returns one parse tree. If there are multiple parse trees for a single string, it prefers `unionl` over `unionr` and also `plus1` over `plus2`. It also never invokes `plus2` with the first string empty, as in this case no progress is made. In fact makes sense to restrict the first string argument of `plus2` to be a `cons-string`—this removes the possibility for a string to have an infinite number of parse trees.

4.5 Parsing

Correctness of `reg2nfa` turns the NFA for a regular expression effectively into a parser. Since we can decide whether a 1×1 matrix contains true or false, using `sound` and `complete`, for any string we can have a parse tree or a proof of that there cannot be one.

$$\text{parse} : (e : \text{RegExp}) \rightarrow (u : \text{String}) \rightarrow u \blacktriangleright e \uplus (u \blacktriangleright e \rightarrow \perp)$$

5 Related Work

Braibant and Pous [1] implement a Coq tactic for deciding equational theory of Kleene algebras. The work is based on checking if two regular expressions represent the same language. This is done in four steps. First, regular expressions are converted into ε -NFAs. Then ε -transitions are removed to get NFAs. Next, determinisation procedure converts NFAs into DFAs. Finally, they check whether the DFAs are equivalent. This results in a general decision procedure for Kleene algebras. In principle, it can be used to solve the recognition problem: to check whether a word w is in the language defined by a regular expression r , we can check if $w \cup r$ and r define the same language. But as this requires go through all four steps for each query, it is impractical.

In contrast, we focus only on the recognition problem. The main difference of our work is that we convert regular expressions directly into NFAs without ε -transitions. This makes the overall process simpler, since then we do not have to find ε -closures and remove ε -transitions afterward and pepper all that, as Braibant and Pous confirm, with quite tricky proofs of correctness.

Many works formalizing recognition of regular languages are based on the concept of the derivative of a language [5,2,3,7]. This is not accidental, since derivatives have nice algebraic properties which make them attractive for a formal development.

It seems that the alternative approach of converting regular expressions to finite automata is believed to be a messy procedure with too much low-level detail involved. For instance, Krauss and Nipkow [5] discuss the link between regular expressions and finite automata in the context of lexing, but point out that encoding finite automata as graphs involves a painful amount of detail and a higher-level approach is desirable.

Wu et al. [9] show how to formalize the Myhill-Nerode theorem by only using regular expressions and the motivation behind this approach is again to avoid the trouble of representing automata as graphs.

We have shown that conversion from regular expressions to finite state automata encoded as Boolean matrices can be done in a concise and high-level way by using block operations on matrices. Proofs in this setting benefit significantly from lemmas about block operations.

6 Conclusion

We presented an implementation of a certified parser generator for regular languages. In particular we showed how to reduce operations and proofs about NFAs into linear-algebra operations and proofs. The practical part of this work was divided into two parts. In the first part, we implemented a generic library for matrices, focusing on block operations. Besides an implementation of basic matrix operations, we also proved many well-known properties of these functions.

In the second part of the practical work, we implemented a transformation of regular expressions to NFAs and proved its correctness. A string is parsed by checking whether the NFA accepts it, as soundness turns the positive answer into a parse tree while completeness can be used to conclude impossibility of a parse tree in the negative case.

This work could be continued in several directions. The implemented framework (the matrix library, RegExp and NFA libraries) can be used to formalize different aspects of regular language theory: minimizing NFAs, showing equivalence of regular expressions, conversion of NFAs to regular expressions, etc.

One variation on the theme of this work would be to consider matrices over natural numbers instead of Booleans. This would allow counting of accepting runs of an NFA (paths from an initial to a final state). Soundness and completeness of the transformation of regular expressions to NFAs would establish a bijection between the parse trees of a given string and the accepting runs of the NFA.

Acknowledgements. This work was supported by the ERDF funded CoE project EXCS, the Estonian Ministry of Education and Research target-financed theme no. 0140007s12 and the Estonian Science Foundation grant no. 9475.

References

1. Braibant, T., Pous, D.: An efficient Coq tactic for deciding Kleene algebras. In: Kaufmann, M., Paulson, L.C. (eds.) ITP 2010. LNCS, vol. 6172, pp. 163–178. Springer, Heidelberg (2010)
2. Coquand, T., Siles, V.: A decision procedure for regular expression equivalence in type theory. In: Jouannaud, J.-P., Shao, Z. (eds.) CPP 2011. LNCS, vol. 7086, pp. 119–134. Springer, Heidelberg (2011)
3. Danielsson, N.A.: Total parser combinators. In: Proc. of 15th ACM SIGPLAN Int. Conf. on Functional Programming, ICFP 2010, pp. 285–296. ACM (2010)
4. Grune, D.: Parsing Techniques: A Practical Guide, 2nd edn. Springer (2010)
5. Krauss, A., Nipkow, T.: Proof pearl: Regular expression equivalence and relation algebra. *J. of Autom. Reasoning* 49(1), 95–106 (2012)
6. Macedo, H.D., Oliveira, J.N.: Typing linear algebra: A biproduct-oriented approach. *Sci. of Comput. Program.* 78(11), 2160–2191 (2013)
7. Morrisett, G., Tan, G., Tassarotti, J., Tristan, J.B., Gan, E.: Rocksalt: better, faster, stronger SFI for the x86. In: Proc. of 33rd ACM SIGPLAN Conf. on Programming Language Design and Implementation, PLDI 2012, pp. 395–404. ACM (2012)
8. Norell, U.: Dependently typed programming in Agda. In: Koopman, P., Plasmeijer, R., Swierstra, D. (eds.) AFP 2008. LNCS, vol. 5832, pp. 230–266. Springer, Heidelberg (2009)
9. Wu, C., Zhang, X., Urban, C.: A formalisation of the Myhill-Nerode theorem based on regular expressions (Proof pearl). In: van Eekelen, M., Geuvers, H., Schmaltz, J., Wiedijk, F. (eds.) ITP 2011. LNCS, vol. 6898, pp. 341–356. Springer, Heidelberg (2011)