

# Validated Parsing of Regular Expressions in Agda

Wai Tak, Cheung  
Student ID: 1465388  
Supervisor: Dr. Martín Escardó



Submitted in conformity with the requirements  
for the degree of BSc. Computer Science  
School of Computer Science  
University of Birmingham

Copyright © 2016 School of Computer Science, University of Birmingham

# Abstract

## Validated Parsing of Regular Expressions in Agda

Wai Tak, Cheung

---

Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah  
blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah  
blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah  
blah. Blah blah blah. Blah blah blah. Blah blah blah.

Keywords: regular expression, finite automata, agda, proof assistant

## Acknowledgments

---

Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah  
blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah  
blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah  
blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah  
blah. Blah blah blah. Blah blah blah. Blah blah blah. Blah blah blah.

Software of this project can be found in  
<https://codex.cs.bham.ac.uk/svn/projects/2015/wtc488/>

## List of Abbreviations

<b><math>\epsilon</math>-NFA</b>	Non-deterministic Finite Automaton with $\epsilon$ -transition
<b>NFA</b>	Non-deterministic Finite Automaton
<b>DFA</b>	Deterministic Finite Automaton
<b>MDFA</b>	Minimised Deterministic Finite Automaton

# Contents

<b>List of Abbreviations</b>	<b>5</b>
<b>1 Introduction</b>	<b>8</b>
1.1 Motivation . . . . .	8
1.2 Overview . . . . .	8
<b>2 Agda</b>	<b>9</b>
2.1 Simply Typed Functional Programming . . . . .	9
2.2 Dependent Types . . . . .	10
2.3 Propositions as Types . . . . .	11
2.3.1 Propositional Logic . . . . .	11
2.3.2 Predicate Logic . . . . .	13
2.3.3 Decidability . . . . .	14
2.3.4 Propositional Equality . . . . .	14
2.4 Program Specifications as Types . . . . .	15
<b>3 Related Work</b>	<b>17</b>
<b>4 Formalisation in Type Theory</b>	<b>18</b>
4.1 Subsets and Decidable Subsets . . . . .	18
4.1.1 Operations on Subsets . . . . .	19
4.1.2 Operations on Decidable Subsets . . . . .	19
4.2 Languages . . . . .	19
4.2.1 Operations on Languages . . . . .	19
4.3 Regular Languages and Regular Expressions . . . . .	20
4.4 $\epsilon$ -Non-deterministic Finite Automata . . . . .	21
4.5 Thompson's Construction . . . . .	24
4.6 Non-deterministic Finite Automata . . . . .	28
4.7 Removing $\epsilon$ -transitions . . . . .	29
4.8 Deterministic Finite Automata . . . . .	29
4.9 Powerset Construction . . . . .	29
4.10 Minimal DFA . . . . .	30
4.11 Minimising DFA . . . . .	30
<b>5 Further Extensions</b>	<b>31</b>

<b>6</b>	<b>Evaluation</b>	<b>32</b>
6.1	Correctness and Readability . . . . .	32
6.2	Different choices of representations . . . . .	33
6.3	Problems arise . . . . .	33
<b>7</b>	<b>Conclusion</b>	<b>35</b>
	<b>Bibliography</b>	<b>36</b>
	<b>Appendices</b>	<b>37</b>

# 1 Introduction

This project aims to study the feasibility of formalising Automata Theory [2] in Type Theory [11] with the aid of a dependently-typed functional programming language, Agda [15]. Automata Theory is an extensive work; therefore, it will be unrealistic to include all the materials under the time constraints. Accordingly, this project will only focus on the theorems and proofs that are related to the translation of regular expressions to finite automata. In addition, this project also serves as an example of how complex and non-trivial proofs are formalised.

Our Agda formalisation is consist of two components: 1) the translation of regular expressions to DFA and 2) the correctness proofs of the translation. At this stage, we are only interested in the correctness of the translation but not the efficiency of the algorithms.

## 1.1 Motivation

My motivation on this project is to learn and apply dependent types in formalising programming logic. At the beginning of this project, I had no knowledge about dependent types. Therefore, we had to choose carefully what theorems to formalise. On one hand, the theorems should be non-trivial enough such that they require a substantial amount of work to be done. On the other hand, they should not be too difficult because I am only a beginner in this area. Finally, we decided to go with the Automata Theory because it satisfies the two criteria and I am familiar with the theory as the basic concepts of the theory was explained in the course *Model of Computation*.

## 1.2 Overview

Section 2 will be a brief introduction on Agda and dependent types. We will describe how Agda can be used as a proof assistant by formalising several small proofs in it. Experienced Agda users can skip this section and start from section 3 directly. In section 3, we will describe several researches that are also related to the formalisation of Automata Theory. Following the background, section 4 will be a detail description of our work. We will walk through the two components in our Agda formalisation. Note that the definitions, theorems and proofs written in this section are extracted from our Agda code. They may be different from their usual mathematical forms in order to adapt the environment in Type Theory. In section 5, we will discuss two possible extensions to our project: 1) Myhill-Nerode Theorem and 2) the Pumping Lemma. After that, in section 6, we will evaluate the project as a whole. Finally, the conclusions will be drawn.



## 2 Agda

Agda is a dependently-typed functional programming language and a proof assistant based on Intuitionistic Type Theory [11]. The current version (Agda 2) is rewritten by Norell [13] during his doctorate study at the Chalmers University of Technology. In this section, we will describe the basic features of Agda and how dependent types are employed to construct programs and proofs. Most of the materials presented below can also be found in the two tutorial papers [4] and [14]. Interested readers can read the two papers in order to get a more precise idea on how to work with Agda. Now, we will begin by showing how to do ordinary functional programming in Agda.

### 2.1 Simply Typed Functional Programming

Haskell is the implementation language of Agda, as shown below, Agda has borrowed many features from Haskell. In the following paragraphs, we will demonstrate how to define basic data types and functions.

**Boolean** We first declare the type of Boolean values in Agda.

```
data Bool : Set where
  true  : Bool
  false : Bool
```

*Bool* has two constructors: *true* and *false*. These two constructors are also elements of *Bool* as they take no arguments. On the other hand, *Bool* itself is a member of the type *Set*. The type of *Set* is *Set*<sub>1</sub> and the type of *Set*<sub>1</sub> is *Set*<sub>2</sub>. The type hierarchy goes on and becomes infinite. Now, let us define the negation of Boolean values.

```
not : Bool → Bool
not true  = false
not false = true
```

Unlike in Haskell, we must provide a type signature explicitly for every function and pattern match on all possible cases of a function in Agda. For instance, the function below will be rejected by the Agda compiler as the case (*not false*) is missing.

```
not : Bool → Bool
not true  = false
```

**Natural Number** Now, let us declare the type of natural numbers in Peano style.

```
data N : Set where
```

```

zero : ℕ
suc   : ℕ → ℕ

```

The constructor **suc** represents the successor of a given natural number. For instance, the number **1** is equivalent to (**suc zero**). Now, let us define the addition of natural numbers recursively as follow:

```

_+_ : ℕ → ℕ → ℕ: Set where
zero + m = m
(suc n) + m = suc (n + m)

```

**Parameterised Types** In Haskell, the type of list **[a]** is parameterised by the type parameter **a**. The analogous data type in Agda is defined as follow:

```

data List (A : Set) : Set where
[]       : List A
_::_     : A → List A → List A

```

Let us try to define a function which takes a list as the argument and returns the first element of the list.

```

head : {A : Set} → List A → A
head [] = {!!}
head (x :: xs) = x

```

What should we return in case **[]**? In Haskell, we can simply skip the **[]** case and the compiler will generate an error for us. However, as mentioned before, we have to pattern match on all possible cases in Agda. One possible workaround here is to return **nothing** of the **Maybe** type for case **[]**. Another solution is to constrain the arguments using dependent types such that the input list will always have at least one element.

## 2.2 Dependent Types

A dependent type is a type that depends on values of other types. For example, **A<sup>n</sup>** is a vector that contains **n** elements of **A**. These kind of types is not possible to be declared in simply-typed systems like Haskell<sup>1</sup> and Ocaml. Now, let us look at how it is declared using dependent type.

```

data Vec (A : Set) : ℕ → Set where
[]       : Vec A zero
_::_     : ∀ {n} → A → Vec A n → Vec A (suc n)

```

---

<sup>1</sup>Haskell itself does not support dependent types by its own. However, there are several APIs in Haskell that simulates dependent types, for example, Ivor [10] and GADT.

In the type signature,  $(A : \mathbf{Set})$  is the type parameter while  $\mathbb{N} \rightarrow \mathbf{Set}$  means that  $\mathbf{Vec}$  takes a number  $n$  from  $\mathbb{N}$  and produces a type that depends on  $n$ . With different natural numbers, we get different types from the inductive family  $\mathbf{Vec}$ . For example,  $\mathbf{Vec} A \mathbf{zero}$  is the type of empty vectors and  $\mathbf{Vec} A \mathbf{10}$  is another vector type with length ten.

Dependent types allow us to be more expressive and precise over type declaration. Let us declare the *head* function for  $\mathbf{Vec}$ .

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

We only need to pattern match on the case  $(x :: xs)$  because the type  $\mathbf{Vec} A (\text{suc } n)$  ensures that the argument will never be  $[]$ . Apart from vectors, we can also define the type of binary search tree in which every tree of this type is guaranteed to be sorted. However, we will not be looking into that in here as it is not our major concern. Interested readers can take a look at Section 6 in [4]. Furthermore, dependent types also allow us to encode predicate logic and program specifications as types. However, before we go into these two applications, we will have to discuss the idea of propositions as types.

## 2.3 Propositions as Types

In the 1930s, Curry identified the correspondence between propositions in propositional logic and types [5]. After that, in the 1960s, de Bruijn and Howard extended Curry's correspondence to predicate logic by introducing dependent types [6, 9]. Later on, Martin-Löf published his work, Intuitionistic Type Theory [11], which turned the correspondence into a new foundational system for constructive mathematics.

In the paragraphs below, we will show how the correspondence is formalised in Agda. Note that Intuitionistic Type Theory is based on constructive logic but not classical logic and there is a fundamental difference between them. Interested readers can take a look at [3]. Now, we will begin with propositional logic.

### 2.3.1 Propositional Logic

In general, Curry's correspondence states that a proposition can be interpreted as a set of its proofs. A proposition is true if and only if its set of proofs is inhabited, i.e. there is at least one element in the set; it is false if and only if its set of proofs is empty.

**Truth** For a proposition to be always true, its corresponding type must have at least one element.

$$\begin{aligned} \text{data } \top &: \mathbf{Set} \text{ where} \\ \text{tt} &: \top \end{aligned}$$

**Falsehood** The proposition that is always false corresponds to a type having no elements at all.

$\text{data } \perp : \text{Set} \text{ where}$

**Conjunction** Suppose  $\mathbf{A}$  and  $\mathbf{B}$  are propositions, then the proofs of their conjunction  $\mathbf{A} \wedge \mathbf{B}$  should contain both a proof of  $\mathbf{A}$  and a proof of  $\mathbf{B}$ . In Type Theory, it corresponds to the product type.

$\text{data } \_ \times \_ (A B : \text{Set}) : \text{Set} \text{ where}$   
 $\_,\_ : A \rightarrow B \rightarrow A \times B$

The above construction resembles the introduction rule of conjunction while the elimination rules are formalised as follow:

$\text{fst} : \{A B : \text{Set}\} \rightarrow A \times B \rightarrow A$   
 $\text{fst } (a, b) = a$

$\text{snd} : \{A B : \text{Set}\} \rightarrow A \times B \rightarrow B$   
 $\text{snd } (a, b) = b$

**Disjunction** Suppose  $\mathbf{A}$  and  $\mathbf{B}$  are propositions, then the proofs of their disjunction  $\mathbf{A} \vee \mathbf{B}$  should contains either a proof of  $\mathbf{A}$  or a proof of  $\mathbf{B}$ . In Type Theory, it is represented by the sum type.

$\text{data } \_ \uplus \_ (A B : \text{Set}) : \text{Set} \text{ where}$   
 $\text{inj}_1 : A \rightarrow A \uplus B$   
 $\text{inj}_2 : B \rightarrow A \uplus B$

The elimination rule of disjunction is defined as follow:

$\uplus\text{-elim} : \{A B C : \text{Set}\}$   
 $\rightarrow A \uplus B$   
 $\rightarrow (A \rightarrow C)$   
 $\rightarrow (B \rightarrow C)$   
 $\rightarrow C$   
 $\uplus\text{-elim } (\text{inj}_1 a) f g = f a$   
 $\uplus\text{-elim } (\text{inj}_2 b) f g = g b$

**Negation** Suppose  $\mathbf{A}$  is a proposition, then its negation is defined as a function that transforms any arbitrary proof of  $\mathbf{A}$  to the falsehood ( $\perp$ ).

$$\neg : \text{Set} \rightarrow \text{Set}$$

$$\neg A = A \rightarrow \perp$$

**Implication** We say that  $A$  implies  $B$  if and only if we can transform every proof of  $A$  into a proof of  $B$ . In Type Theory, it corresponds to a function from  $A$  to  $B$ , i.e.  $A \rightarrow B$ .

**Equivalence** Two propositions  $A$  and  $B$  are equivalent if and only if  $A$  implies  $B$  and  $B$  implies  $A$ . We can consider it as a conjunction of the two implications.

$$\_ \Longleftrightarrow \_ : \text{Set} \rightarrow \text{Set} \rightarrow \text{Set}$$

$$A \Longleftrightarrow B = (A \rightarrow B) \times (B \rightarrow A)$$

Now, by using the above constructions, we can formalise theorems of propositional logic in Agda. For example, we can prove that if  $P$  implies  $Q$  and  $Q$  implies  $R$ , then  $P$  implies  $R$ .

$$\begin{aligned} \text{prop-lem} & : \{P\ Q : \text{Set}\} \\ & \rightarrow (P \rightarrow Q) \\ & \rightarrow (Q \rightarrow R) \\ & \rightarrow (P \rightarrow R) \\ \text{prop-lem } f\ g & = \lambda p \rightarrow g\ (f\ p) \end{aligned}$$

By completing the function, we have provided an element to the type  $(P \rightarrow Q) \rightarrow (Q \rightarrow R) \rightarrow (P \rightarrow R)$  and thus, we have also proved the theorem to be true.

### 2.3.2 Predicate Logic

We will now move on to predicate logic and introduce the universal ( $\forall$ ) and existential ( $\exists$ ) quantifiers. A predicate is represented by a dependent type in the form of  $A \rightarrow \text{Set}$ . For example, we can define the predicate of even numbers and odd numbers inductively as follow:

$$\begin{aligned} & \text{mutual} \\ & \text{data } \_ \text{isEven} : \mathbb{N} \rightarrow \text{Set} \text{ where} \\ & \quad \text{base} : \text{zero isEven} \\ & \quad \text{step} : \forall n \rightarrow n \text{ isOdd} \rightarrow (\text{suc } n) \text{ isEven} \\ & \text{data } \_ \text{isOdd} : \mathbb{N} \rightarrow \text{Set} \text{ where} \\ & \quad \text{step} : \forall n \rightarrow n \text{ isEven} \rightarrow (\text{suc } n) \text{ isOdd} \end{aligned}$$

**Universal Quantifier** The interpretation of the universal quantifier is similar to implication. In order for  $\forall x \in A. B(x)$  to be true, we will have to transform every proof  $(a)$  of  $A$  into a proof of the predicate  $B[x := a]$ . In Type Theory, it is represented by the function  $(x : A) \rightarrow B\ x$ . For example, we can prove by induction that for every natural number, it is either even or odd.

```

lem1 : ∀ n → n isEven ∨ n isOdd
lem1 zero = inj1 base
lem1 (suc n) with lem1 n
... | inj1 nIsEven = inj2 (step n nIsEven)
... | inj2 nIsOdd = inj1 (step n nIsOdd)

```

**Existential Quantifier** The interpretation of the existential quantifier is similar to conjunction. In order for  $\exists x \in A. B(x)$  to be true, we need to provide a proof  $(a)$  of  $A$  and a proof  $(p)$  of the predicate  $B[x := a]$ . In Type Theory, it is represented by the generalised product type  $\Sigma$ .

```

data Σ (A : Set) (B : A → Set) : where
  _,_ : (a : A) → B a → Σ A B

```

For simplicity, we will change the syntax of  $\Sigma$  to  $\exists[x \in A] B$ . As an example, we can prove that there exists a natural number which is even.

```

lem2 : ∃[ n ∈ ℕ ] (n isEven)
lem2 = zero , base

```

### 2.3.3 Decidability

A proposition  $P$  is decidable if and only if there exists an algorithm that can decide whether it is true or false. We will define the decidability of a proposition as follow:

```

data Dec (A : Set) : Set where
  yes : A → Dec A
  no  : ¬ A → Dec A

```

### 2.3.4 Propositional Equality

One important feature of Type Theory is that we can define the equality of propositions as types. The equality relation is interpreted as follow:

```

data _≡_ {A : Set} (x : A) : A → Set where
  refl : x ≡ x

```

This states that for any  $x$  in  $\mathbf{A}$ ,  $\mathbf{refl}$  is an element of the type  $x \equiv x$ . More generally,  $\mathbf{refl}$  is a proof of  $x \equiv x'$  provided that  $x$  and  $x'$  is the same after normalisation. For example, we can prove that  $\exists n \in \mathbb{N}. n = 1 + 1$  as follow:

$$\begin{aligned} \text{lem}_3 & : \exists[ n \in \mathbb{N} ] n \equiv (1 + 1) \\ \text{lem}_3 & = \text{suc } (\text{suc } \text{zero}) , \text{refl} \end{aligned}$$

We can put  $\mathbf{refl}$  in the proof only because both  $\mathbf{suc } (\mathbf{suc } \text{zero})$  and  $1 + 1$  have the same form after normalisation. Now, let us define the elimination rule of equality. The rule should allow us to substitute equivalence objects into any proposition.

$$\begin{aligned} \text{subst} & : \{A : \text{Set}\} \{x y : A\} \rightarrow (P : A \rightarrow \text{Set}) \rightarrow x \equiv y \rightarrow P x \rightarrow P y \\ \text{subst } P \text{ refl } p & = p \end{aligned}$$

We can also prove the congruency of equality:

$$\begin{aligned} \text{cong} & : \{A B : \text{Set}\} \{x y : A\} \rightarrow (f : A \rightarrow B) \rightarrow x \equiv y \rightarrow f x \equiv f y \\ \text{cong } f \text{ refl} & = \text{refl} \end{aligned}$$

## 2.4 Program Specifications as Types

As we mentioned before, dependent types also allow us to encode program specifications within the same platform. In order to demonstrate the idea, we will give an example that uses sorted lists of natural numbers. Let us begin by defining a predicate of sorted list (in ascending order).

$$\begin{aligned} \text{All-lt} & : \mathbb{N} \rightarrow \text{List } \mathbb{N} \rightarrow \text{Set} \\ \text{All-lt } n [] & = \top \\ \text{All-lt } n (x :: xs) & = n \leq x \times \text{All-lt } n xs \end{aligned}$$

$$\begin{aligned} \text{Sorted-ASC} & : \text{List } \mathbb{N} \rightarrow \text{Set} \\ \text{Sorted-ASC} [] & = \top \\ \text{Sorted-ASC } (x :: xs) & = \text{All-lt } x xs \times \text{Sorted-ASC } xs \end{aligned}$$

$\text{All-lt}$  defines the condition where a given number is smaller than all the numbers inside a given list. Now, let us define an insertion function that takes a natural number and a list as the arguments and returns a list of natural numbers. The insertion function is designed in a way that if the input list is already sorted, then the output list will also be sorted.

$$\begin{aligned} \text{insert} & : \mathbb{N} \rightarrow \text{List } \mathbb{N} \rightarrow \text{List } \mathbb{N} \\ \text{insert } n [] & = n :: [] \\ \text{insert } n (x :: xs) & \text{ with } n \leq? x \\ \dots \mid \text{yes } \_ & = n :: (x :: xs) \\ \dots \mid \text{no } \_ & = x :: \text{insert } n xs \end{aligned}$$

Note that  $\_ \leq? \_$  has the type  $\forall \mathbf{n} \mathbf{m} \rightarrow \mathbf{Dec} (\mathbf{n} \leq \mathbf{m})$ . It is a proof of the decidability of  $\_ \leq \_$  and it can also be used to determine whether a given number  $\mathbf{n}$  is less than or equal to another number  $\mathbf{m}$ . Now, let us encode the specification of the insertion function as follow:

$$\begin{aligned} \text{insert-sorted} &: \forall \{n\} \{as\} \\ &\rightarrow \text{Sorted-ASC } as \\ &\rightarrow \text{Sorted-ASC } (\text{insert } n \text{ as}) \end{aligned}$$

In the type signature,  $(\text{Sorted-ASC } as)$  corresponds to the pre-condition and  $(\text{Sorted-ASC } (\text{insert } n \text{ as}))$  corresponds to the post-condition. Once we have completed the function, we will also have proved the specification to be true. Interested readers are recommended to finished the proof.



### 3 Related Work

Agular and Mannaa published a similar work [1] in 2009. They constructed a decider for regular expressions which can determine whether a given string is accepted by a given regular expression. The decider was based on the calculation of the derivation of a regular expression which only needs to convert the regular expression into part of an automaton. Their decider was implemented using the *Maybe* type as follow:

$$\text{accept} : (\text{re} : \text{RegExp}) \rightarrow (\text{as} : \text{List carrier}) \rightarrow \text{Maybe } (\text{as} \in \sim \llbracket \text{re} \rrbracket)$$

When a string is accepted by the regular expression, i.e.  $w \in L(e)$ , the decider will return its proof. However it fails to generate a proof for the opposite case, i.e.  $w \notin L(e)$ . As they explained in the paper, it is not possible without converting the regular expression into the entire finite automaton.

While in 2013, Firsov and Uustalu also published another related research paper [7]. They translated regular expressions into NFA and proved that their accepting languages are equal. Unlike Agular and Mannaa's decider, Firsov and Uustalu's algorithm could generate proofs for both cases when a string is or is not accepted by the regular expression. In their definition of NFA, the set of states  $Q$  and its subsets are represented as vectors while the transition function  $\delta$  takes an alphabet as the argument and returns a matrix representation of the transition table.

$$\begin{aligned} \text{record NFA} &: \text{Set where} \\ \text{field} & \\ |Q| &: \mathbb{N} \\ \delta &: \Sigma \rightarrow |Q| * |Q| \\ I &: 1 * |Q| \\ F &: |Q| * 1 \end{aligned}$$

Note that  $_ * _$  is an inductive family that takes two natural numbers  $n$  and  $m$  and produces a matrix type  $n \times m$ . This representation allows us to iterate the set easily but it looks unnatural compare to the actual mathematical definition of NFA.

## 4 Formalisation in Type Theory

Let us recall the two components of the formalisation: 1) translating any regular expressions to a DFA and 2) proving the correctness of the translation.

In part 1), the translation was divided into the following steps. First, we followed Thompson's construction algorithm to convert any regular expressions to an  $\epsilon$ -NFA. Then we removed all the  $\epsilon$ -transitions in the  $\epsilon$ -NFA by computing the  $\epsilon$ -closure for every states. After that, we used powerset construction to create a DFA. Finally, we removed all the unreachable states and then used quotient construction to obtain the minimised DFA.

In part 2), the correctness proofs of the above translation were also separated into different steps according to part 1). For each of the translation steps in part 1), we proved that the language accepted by the input is equal to the language accepted by its translated output. i.e.  $L(regex) = L(translated \ \epsilon\text{-NFA}) = L(translated \text{ DFA}) = L(translated \text{ MDFA})$ .

In the following parts, we will walk through the formalisation of each of the above steps together with their correctness proofs. Note that all the definitions, theorems, lemmas and proofs written in below are adapted to the formalisation in Agda. Now, before we go into regular expressions and automata, we first need to have a representation of subsets and languages as they are fundamental elements in the theory.

### 4.1 Subsets and Decidable Subsets

**Definition 1.1** Suppose  $A$  is a set, in Type Theory, its subsets are represented as a unary function on  $A$ , i.e.  $Subset \ A = A \rightarrow Set$ .

When declaring a subset in Agda, we can write  $SubA = \lambda a \rightarrow ?$ , the  $?$  here defines the condition for  $a$  to be included in  $SubA$ . This construction is very similar to set comprehension. For example, the subset  $\{a \mid a \in A, P(a)\}$  corresponds to  $\lambda a \rightarrow P \ a$ . Subset is also a unary predicate of  $A$ ; therefore, the decidability of it will remain unknown until it is proved.

**Definition 1.2** The other representation of subset is  $DecSubset \ A = A \rightarrow Bool$ . Unlike  $Subset$ , its decidability is ensured by its definition.

The two definitions have different purposes.  $Subset$  is used to represent *Language* because not every language is decidable. For other parts such as a subset of states in an automaton,  $DecSubset$  is used as the decidability is assumed in the definition. The two definitions are defined in `Subset.agda`

and `Subset/DecidableSubset.agda` respectively as stated at the top. Operators such as membership ( $\in$ ), subset ( $\subseteq$ ), superset ( $\supseteq$ ) and equality ( $=$ ) can also be found in the two files.

#### 4.1.1 Operations on Subsets

#### 4.1.2 Operations on Decidable Subsets

Now, by using the representation of subset, we can define languages, regular expressions and finite automata.

### 4.2 Languages

Suppose we have a set of alphabets  $\Sigma$ ; in Type Theory, it can be represented as a data type, i.e.  $\Sigma : \text{Set}$ . Notice that the decidable equality of  $\Sigma$  is assumed. In Agda, they are passed to every modules as parameters  $(\Sigma : \text{Set}) (dec : \text{DecEq } \Sigma)$ .

**Definition 2.1** We first define  $\Sigma^*$  as the set of all strings over  $\Sigma$ . In our approach, it was expressed as a list of  $\Sigma$ , i.e.  $\Sigma^* = \text{List } \Sigma$ .

For example,  $(A :: g :: d :: a :: [])$  represents the string 'Agda' and the empty list  $[]$  represents the empty string  $\epsilon$ . In this way, we can pattern match on the input string in order to get the first input alphabet and to run a transition from a particular state to another state.

**Definition 2.2** A language is a subset of  $\Sigma^*$ ; in Type Theory,  $\text{Language} = \text{Subset } \Sigma^*$ . Notice that *Subset* instead of *DecSubset* is used because not every language is decidable.

#### 4.2.1 Operations on Languages

**Definition 2.3** If  $L_1$  and  $L_2$  are languages, then the union of the two languages  $L_1 \cup L_2$  is defined as  $\{w \mid w \in L_1 \vee w \in L_2\}$ . In Type Theory, we define it as  $L_1 \cup L_2 = \lambda w \rightarrow w \in L_1 \uplus w \in L_2$ .

**Definition 2.4** If  $L_1$  and  $L_2$  are languages, then the concatenation of the two languages  $L_1 \bullet L_2$  is defined as  $\{w \mid \exists u \in L_1. \exists v \in L_2. w = uv\}$ . In Type Theory, we define it as  $L_1 \bullet L_2 = \lambda w \rightarrow \exists [u \in \Sigma^*] \exists [v \in \Sigma^*] (u \in L_1 \times v \in L_2 \times w \equiv u ++ v)$ .

**Definition 2.5** If  $L$  is a language, then the closure of  $L$ ,  $L^*$  is defined as  $\bigcup_{n \in \mathbb{N}} L^n$  where  $L^n = L \bullet L^{n-1}$  and  $L^0 = \{\epsilon\}$ . In Type Theory, we have  $L \star = \lambda w \rightarrow \exists[n \in \mathbb{N}](w \in L \wedge n)$  where the function  $\wedge$  is defined recursively as:

$$\begin{aligned} \wedge & : \text{Language} \rightarrow \text{Language} \rightarrow \text{Language} \\ L \wedge \text{zero} & = \llbracket \epsilon \rrbracket \\ L \wedge (\text{suc } n) & = L \bullet L \wedge n \end{aligned}$$

### 4.3 Regular Languages and Regular Expressions

**Definition 3.1** We define regular languages over  $\Sigma$  inductively as follow:

1.  $\emptyset$  is a regular language;
2.  $\{\epsilon\}$  is a regular language;
3.  $\forall a \in \Sigma$ .  $\{a\}$  is a regular language;
4. if  $L_1$  and  $L_2$  are regular languages, then
  - (a)  $L_1 \cup L_2$  is a regular language;
  - (b)  $L_1 \bullet L_2$  is a regular language;
  - (c)  $L_1 \star$  is a regular language.

Listing 1: Regular languages

```
data Regular : Language → Set1 where
  nullL : ∀ {L} → L ≈ ∅ → Regular L
  empty : ∀ {L} → L ≈ ⟦ε⟧ → Regular L
  singl : ∀ {L} → (a : Σ) → L ≈ ⟦a⟧ → Regular L
  union : ∀ {L} L1 L2 → Regular L1 → Regular L2 → L ≈ L1 ∪ L2 → Regular L
  conca : ∀ {L} L1 L2 → Regular L1 → Regular L2 → L ≈ L1 • L2 → Regular L
  kleen : ∀ {L} L1 → Regular L1 → L ≈ L1 ★ → Regular L
```

**Definition 3.2** Here we define regular expressions inductively over  $\Sigma$  as follow:

1.  $\emptyset$  is a regular expression denoting the regular language  $\emptyset$ ;
2.  $\epsilon$  is a regular expression denoting the regular language  $\{\epsilon\}$ ;
3.  $\forall a \in \Sigma$ .  $a$  is a regular expression denoting the regular language  $\{a\}$ ;
4. if  $e_1$  and  $e_2$  are regular expressions denoting the regular languages  $L_1$  and  $L_2$  respectively, then
  - (a)  $e_1 \mid e_2$  is a regular expressions denoting the regular language  $L_1 \cup L_2$ ;
  - (b)  $e_1 \cdot e_2$  is a regular expression denoting the regular language  $L_1 \bullet L_2$ ;
  - (c)  $e_1^*$  is a regular expression denoting the regular language  $L_1 \star$ .

The Agda formalisation is separated into two parts, firstly the definition of regular expressions and secondly the languages denoted by them.

Listing 2: Regular expressions

```
data RegExp : Set where
  ∅      : RegExp
  ε      : RegExp
  σ      : Σ → RegExp
  _|_    : RegExp → RegExp → RegExp
  _·_    : RegExp → RegExp → RegExp
  _*     : RegExp → RegExp
```

Listing 3: Languages denoted by regular expressions

```
LR : RegExp → Language
LR ∅ = ∅
LR ε = [ε]
LR (σ a) = [a]
LR (e1 | e2) = LR e1 ∪ LR e2
LR (e1 · e2) = LR e1 • LR e2
LR (e*) = (LR e) ★
```

#### 4.4 $\epsilon$ -Non-deterministic Finite Automata

By now, the set of strings we have considered are in the form of  $List\ \Sigma^*$ . However, this definition gives us no way to extract an  $\epsilon$ -transition from the input string. Therefore, we need to introduce another representation of the set of strings specifically for this purpose. (For Definition 4.1 and 4.2, please refers to Language.agda)

**Definition 4.1** We define  $\Sigma^\epsilon$  as the union of  $\Sigma$  and  $\{\epsilon\}$ , i.e.  $\Sigma^\epsilon = \Sigma \cup \{\epsilon\}$ .

In Agda, this can be expressed by a data type definition:

```
data Σε : Set where
  α : Σ → Σε
  E : Σε
```

**Definition 4.2** Now we define  $\Sigma^{e*}$ , the set of all strings over  $\Sigma^\epsilon$  in a way similar to  $\Sigma^*$ , i.e.  $\Sigma^{e*} = List\ \Sigma^\epsilon$ .

For example, the string 'Agda' can be represented by  $(\alpha A :: \alpha g :: E :: \alpha d :: E :: \alpha a :: [])$  or  $(E :: \alpha A :: E :: E :: \alpha g :: \alpha d :: E :: \alpha a :: [])$ . We say that these two lists are  $\epsilon$ -strings of the word 'Agda'. When pattern matching on an  $\epsilon$ -string, we can know if there is an  $\epsilon$ -transition or not. Other operators and lemmas regarding  $\epsilon$ -strings such as  $to\Sigma^* : \Sigma^{e*} \rightarrow \Sigma^*$  can also be found in `Language.agda`.

Now, let us define  $\epsilon$ -NFA.

**Definition 4.3** An  $\epsilon$ -NFA is a 5-tuple  $M = (Q, \Sigma^e, \delta, q_0, F)$ , where

1.  $Q$  is a finite set of states;
2.  $\Sigma^e$  is the union of  $\Sigma$  and  $\{\epsilon\}$ ;
3.  $\delta$  is a mapping from  $Q \times \Sigma^e$  to  $\mathcal{P}(Q)$  which defines the behaviour of the automata;
4.  $q_0$  in  $Q$  is the initial state;
5.  $F \subseteq Q$  is the set of accepting states.

Listing 4:  $\epsilon$ -NFA

```
record  $\epsilon$ -NFA : Set1 where
  field
    Q      : Set
     $\delta$     : Q  $\rightarrow$   $\Sigma^e \rightarrow$  DecSubset Q
    q0    : Q
    F      : DecSubset Q
     $\forall qEq$  :  $\forall q \rightarrow q \in^d \delta \ q \ E$ 
    Q?     : DecEq Q
    |Q|-1  :  $\mathbb{N}$ 
    It     : Vec Q (suc |Q|-1)
     $\forall q \in It$  : (q : Q)  $\rightarrow$  (q  $\in^V$  It)
    unique : Unique It
```

The set of alphabets  $\Sigma : Set$  is passed to the file parameters. Together with  $Q$ ,  $\delta$ ,  $q_0$  and  $F$ , these five fields correspond to the 5-tuple  $\epsilon$ -NFA.  $\forall qEq$  is a proof that any state in  $Q$  can reach itself by an  $\epsilon$ -transition.  $Q?$  is the decidable equality of  $Q$ .  $|Q| - 1$  is the number of states - 1. ' $It$ ' is a vector of length  $|Q|$  containing all the states in  $Q$ .  $\forall q \in It$  is a proof that all states in  $Q$  are also in the vector ' $It$ '. *unique* is a proof that there is no repeating elements in ' $It$ '. These extra fields are important when computing  $\epsilon$ -closures, we will look into them again later in more details.

Now, we want to define the set of strings  $\Sigma^*$  accepted by a given  $\epsilon$ -NFA. However, before we can do this, we have to define some operations.

**Definition 4.4** A configuration is a pair  $Q \times \Sigma^{e*}$ . Notice that the configuration is based on  $\Sigma^{e*}$  but not  $\Sigma^*$ .

**Definition 4.5** A move by an  $\epsilon$ -NFA  $N$  is represented by a binary function  $\vdash$  on configurations. We say that  $(q, aw) \vdash (q', w)$  for all  $w$  in  $\Sigma^{e*}$  if and only if  $q' \in \delta(q, a)$  where  $a \in \Sigma^e$ .

$$\begin{aligned} \_ \vdash \_ &: (Q \times \Sigma^e \times \Sigma^{e*}) \rightarrow (Q \times \Sigma^{e*}) \rightarrow \text{Set} \\ (q, a, w) \vdash (q', w') &= w \equiv w' \times q' \in^d \delta(q, a) \end{aligned}$$

**Definition 4.6** We say that  $C \vdash^0 C'$  if and only if  $C = C'$ . We say that  $C_0 \vdash^k C_k$  for any  $k \geq 1$  if and only if there exists a chain of configurations  $C_1, C_2, \dots, C_{k-1}$  such that  $C_i \vdash C_{i+1}$  for all  $0 \leq i \leq k$ .

$$\begin{aligned} \_ \vdash^k \_ &: (Q \times \Sigma^{e*}) \rightarrow \mathbb{N} \rightarrow (Q \times \Sigma^{e*}) \rightarrow \text{Set} \\ (q, w^e) \vdash^k \text{zero} &- (q', w'^e) \\ &= q \equiv q' \times w^e \equiv w'^e \\ (q, w^e) \vdash^k \text{suc } n &- (q', w'^e) \\ &= \exists [p \in Q] \exists [a^e \in \Sigma^e] \exists [u^e \in \Sigma^{e*}] \\ &\quad (w^e \equiv a^e :: u^e \times (q, a^e, u^e) \vdash (p, u^e) \times (p, u^e) \vdash^k n - (q', w'^e)) \end{aligned}$$

**Definition 4.7** We say that  $C \vdash^* C'$  if and only if there exists a number of chains  $n$  such that  $C \vdash^n C'$ .

$$\begin{aligned} \_ \vdash^* \_ &: (Q \times \Sigma^{e*}) \rightarrow (Q \times \Sigma^{e*}) \rightarrow \text{Set} \\ (q, w^e) \vdash^* (q', w'^e) &= \exists [n \in \mathbb{N}] (q, w^e) \vdash^k n - (q', w'^e) \end{aligned}$$

**Definition 4.8** For any string  $w$ , it is accepted by an  $\epsilon$ -NFA  $N$  if and only if there exists a chain of configurations from  $q_0, w^e$  to  $q, \epsilon$  where  $w^e$  is an  $\epsilon$ -string of  $w$  and  $q \in F$ .

**Definition 4.9** The language accepted by an  $\epsilon$ -NFA is given by the set  $\{ w \mid \exists w^e \in \Sigma^{e*}. w = \text{to}\Sigma^*(w^e) \wedge \exists q \in F. (q_0, w^e) \vdash^* (q, \epsilon) \}$ .

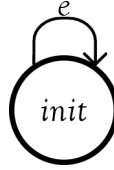
$$\begin{aligned} L^{eN} &: \epsilon\text{-NFA} \rightarrow \text{Language} \\ L^{eN} \text{ nfa} &= \lambda w \rightarrow \\ &\quad \exists [w^e \in \Sigma^{e*}] (w \equiv \text{to}\Sigma^* w^e \times (\exists [q \in Q] (q \in^d F \times (q_0, w^e) \vdash^* (q, [])))) \end{aligned}$$

Now that we have the definition of regular expressions and  $\epsilon$ -NFA, we can formulate the translation using Thompson's Construction.

## 4.5 Thompson's Construction

**Definition 5.1** The translation for any regular expressions to an  $\epsilon$ -NFA is defined inductively as follow:

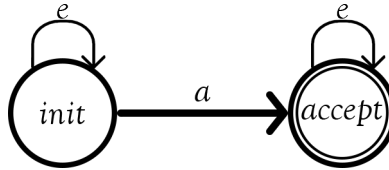
1. for  $\emptyset$ , we have  $M = (\{init\}, \Sigma^e, \delta, init, \emptyset)$  and graphically



2. for  $\epsilon$ , we have  $M = (\{init\}, \Sigma^e, \delta, init, \{init\})$  and graphically



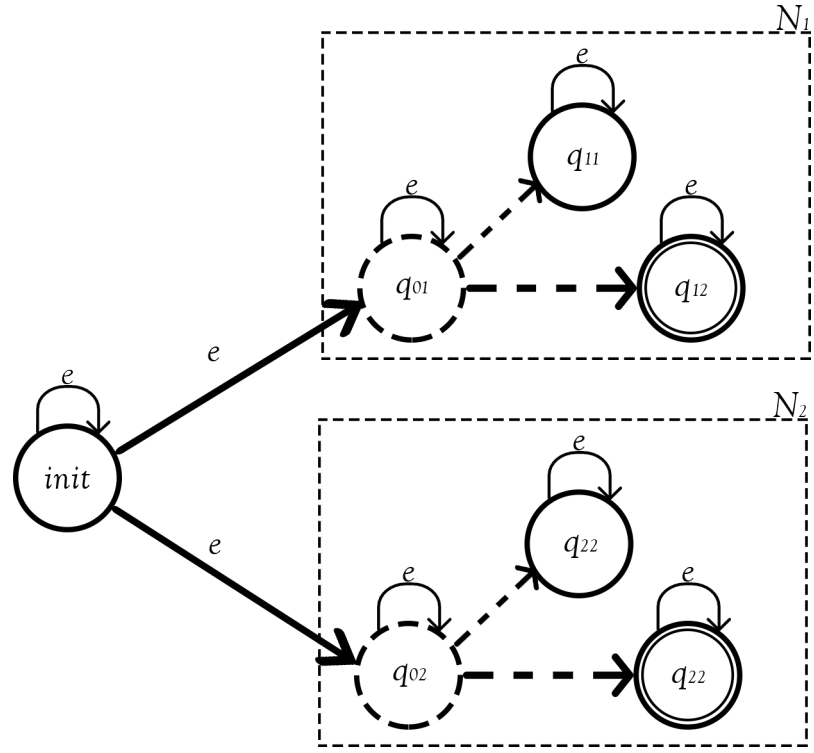
3. for  $a$ , we have  $M = (\{init, accept\}, \Sigma^e, \delta, init, \{accept\})$  and graphically



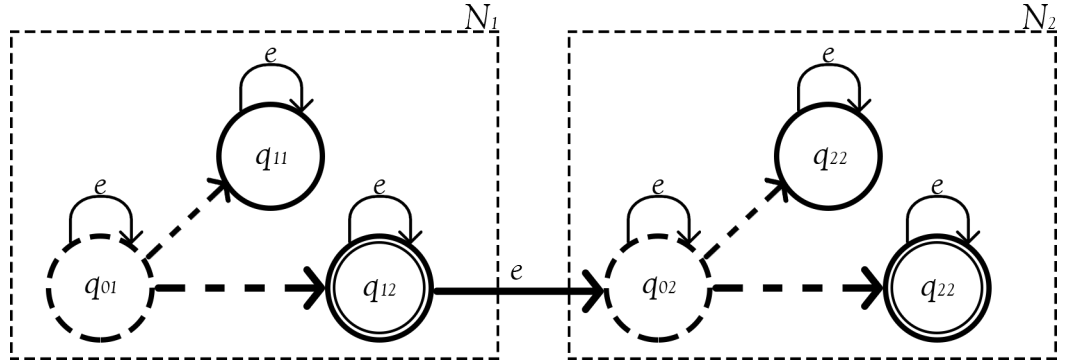
4. if  $N_1 = (Q_1, \delta_1, q_{01}, F_1)$  and  $N_2 = (Q_2, \delta_2, q_{02}, F_2)$  are  $\epsilon$ -NFAs translated from the regular expressions  $e_1$  and  $e_2$  respectively, then

- (a) for  $(e_1 \mid e_2)$ , we have  $M = (\{init\} \cup Q_1 \cup Q_2, \Sigma^e, \delta, init, F_1 \cup F_2)$  and graphically

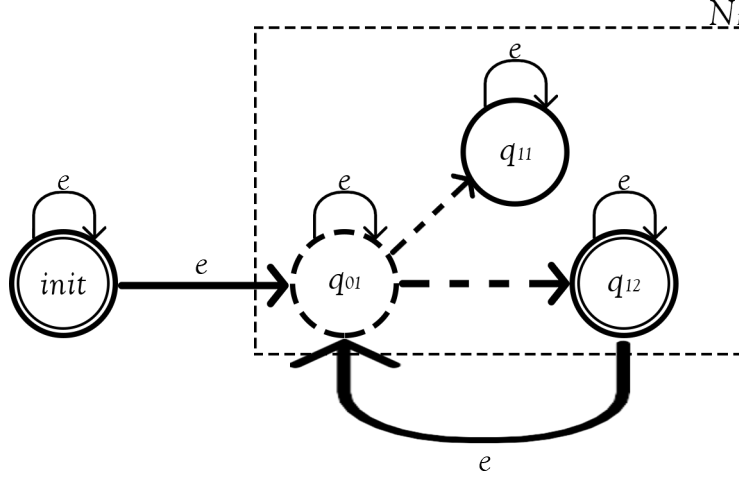




(b) for  $e_1 \cdot e_2$ , we have  $M = (Q_1 \cup \{mid\} \cup Q_2, \Sigma^e, \delta, init, F_2)$  and graphically



(c) for  $e_1^*$ , we have  $M = (\{init\} \cup Q_1, \Sigma^e, \delta, init, \{init\} \cup F_1)$  and graphically



**Theorem 1.1** For any given regular expressions, its accepted language is equal to the language accepted by its translated  $\epsilon$ -NFA using Thompson's Construction. i.e.  $L(e) = L(\text{translated } \epsilon\text{-NFA})$ .

**Proof 1.1** We have to prove that for any regular expressions  $e$ ,  $L(e) \subseteq L(\text{translated } \epsilon\text{-NFA})$  and  $L(e) \supseteq L(\text{translated } \epsilon\text{-NFA})$  by induction on  $e$ .

**Base cases** For  $\emptyset$ ,  $\epsilon$  and  $a$ , by Definition 5.1, it is obvious that the language accepted by them are equal to the language accepted by their translated  $\epsilon$ -NFA.

**Induction hypothesis** For any regular expressions  $e_1$  and  $e_2$ , let  $N_1 = (Q_1, \delta_1, q_{01}, F_1)$  and  $N_2 = (Q_2, \delta_2, q_{02}, F_2)$  be their translated  $\epsilon$ -NFA using Definition 5.1 respectively. Then we assume that  $L(e_1) = L(N_1)$  and  $L(e_2) = L(N_2)$ .

#### Inductive steps

1) For  $(e_1 \mid e_2)$ , let  $M = (Q, \delta, q_0, F) = (\{init\} \cup Q_1 \cup Q_2, \delta, init, F_1 \cup F_2)$  be its translated  $\epsilon$ -NFA using Definition 5.1. Then for any string  $w$ ,

1.1) if  $(e_1 \mid e_2)$  accepts  $w$ , by Definition 3.2, either i)  $e_1$  accepts  $w$  or ii)  $e_2$  accepts  $w$ . Assuming case i), then by induction hypothesis,  $N_1$  also accepts  $w$  which also implies that there exists a chain  $(q_{01}, w^\epsilon) \vdash^* (q, \epsilon)$  in  $N_1$  such that  $w^\epsilon$  is an  $\epsilon$ -string of  $w$  and  $q \in F_1$ . Now, we can add an  $\epsilon$ -transition from  $init$  to  $q_{01}$  in  $M$  such that  $(init, \epsilon w^\epsilon) \vdash^* (q, \epsilon)$  because  $q_{01} \in \delta init \epsilon$ . Now, since  $q \in F_1$  implies that  $q \in F$  and  $\epsilon w^\epsilon$  is also an  $\epsilon$ -string of  $w$ ; therefore  $w \in L(M)$ . The same argument also applies for the case when  $e_2$  accepts  $w$ . Since we have proved that  $w \in L(e_1 \mid e_2) \Rightarrow w \in L(M)$ ; therefore  $L(e_1 \mid e_2) \subseteq L(M)$  also follows;

1.2) if  $M$  accepts  $w$ , then there must exists a chain  $(init, w^\epsilon) \vdash^* (q, \epsilon)$  in  $M$  such that  $w^\epsilon$  is an  $\epsilon$ -string of  $w$  and  $q \in F$ . Since  $q \in F$ , therefore  $q \neq init$ . By Definition 5.1, there are only two possible ways for  $init$  to reach  $q$ , via  $q_{01}$  or ii)  $q_{02}$ . Assuming case i), then we have  $(init, \epsilon^+ w_1) \vdash^* (q_{01}, w_1)$  and  $(q_{01}, w_1) \vdash^* (q, \epsilon)$  where  $w^\epsilon = \epsilon^+ w_1$  and  $q \in Q_1$ . Since we have  $q \in F$  and  $q \in Q_1$ ; therefore we have  $q \in F_1$ . Also  $w_1$  is also an  $\epsilon$ -string of  $w$ , thus the chain  $(q_{01}, w_1) \vdash^* (q, \epsilon)$  implies that  $w \in L(N_1)$ . By induction hypothesis, we have  $w \in L(e_1)$  and thus  $w \in L(e_1 \mid e_2)$ . The same argument also applies for case ii). Since we have proved that  $w \in L(M) \Rightarrow w \in L(e_1 \mid e_2)$ ; therefore  $L(e_1 \mid e_2) \supseteq L(M)$  also follows;

1.3) combining 1.1 and 1.2, we have  $L(e_1 \mid e_2) = L(M)$ .

2) For  $(e_1 \cdot e_2)$ , let  $M = (Q, \delta, q_0, F) = (Q_1 \cup \{mid\} \cup Q_2, \delta, q_{01}, F_2)$  be its translated  $\epsilon$ -NFA using Definition 5.1. Then for any string  $w$ ,

2.1) if  $(e_1 \cdot e_2)$  accepts  $w$ , then by Definition 3.2, there exists a  $u \in L(e_1)$  and a  $v \in L(e_2)$  such that  $w = uv$ . By induction hypothesis,  $u \in L(e_1)$  implies that  $u \in L(N_1)$  and  $v \in L(e_2)$  implies that  $v \in L(N_2)$ . So there exists a chain: i)  $(q_{01}, u^\epsilon) \vdash^* (q_1, \epsilon)$  in  $N_1$  where  $u^\epsilon$  is an  $\epsilon$ -string of  $u$  and  $q_1 \in F_1$  and ii)  $(q_{02}, v^\epsilon) \vdash^* (q_2, \epsilon)$  in  $N_2$  where  $v^\epsilon$  is an  $\epsilon$ -string of  $v$  and  $q_2 \in F_2$ . Now we can add an  $\epsilon$ -transition from  $q_1$  to  $mid$  and from  $mid$  to  $q_{02}$  in order to construct a chain in  $M$ . Since  $q_2 \in F_2$  implies that  $q_2 \in F$  and  $u^\epsilon v^\epsilon$  is an  $\epsilon$ -string of  $w$  implies that so is  $u^\epsilon \epsilon v^\epsilon$ ; therefore  $w \in L(M)$ . Since we have proved that  $w \in L(e_1 \cdot e_2) \Rightarrow w \in L(M)$ , therefore  $L(e_1 \cdot e_2) \subseteq L(M)$  also follows;

2.2) if  $M$  accepts  $w$ , then by Definition 5.1, there must exists a chain  $(init, w^\epsilon) \vdash^* (q, \epsilon)$  in  $M$  where  $w^\epsilon$  is an  $\epsilon$ -string of  $w$  and  $q \in F$ . Since  $q \in F$ , so  $q$  must also be in  $Q_2$ . The only possible way for  $q_{01}$  to reach  $q$  is to go through  $mid$ . This implies that there exists a  $q_1 \in Q_1$ , a  $u^\epsilon \in \Sigma^{\epsilon*}$  and a  $v^\epsilon \in \Sigma^{\epsilon*}$  such that  $(q_{01}, u^\epsilon \epsilon^+ \epsilon^+ v^\epsilon) \vdash^* (q_1, \epsilon^+ \epsilon^+ v^\epsilon)$ ,  $q_1 \in F_1$ ,  $(q_{02}, v^\epsilon) \vdash^* (q_2, \epsilon)$  and  $w^\epsilon = u^\epsilon \epsilon^+ \epsilon^+ v^\epsilon$ . Let  $u$  and  $v$  be the strings represented by  $u^\epsilon$  and  $v^\epsilon$  respectively, we have  $u \in L(N_1)$  and  $v \in L(N_2)$ . Then, by induction hypothesis,  $u \in L(e_1)$  and  $v \in L(e_2)$ . Since  $w^\epsilon$  is an  $\epsilon$ -string of  $w$ , so is  $u^\epsilon v^\epsilon$  and thus  $w = uv$ . From this, we can deduce that  $w \in L(e_1 \cdot e_2)$ . Since we have proved that  $w \in L(M) \Rightarrow w \in L(e_1 \cdot e_2)$ , therefore  $L(e_1 \cdot e_2) \supseteq L(M)$  also follows;

2.3) combining 2.1 and 2.2, we have  $L(e_1 \cdot e_2) = L(M)$ .

3) For  $e^*$ , let  $M = (Q, \delta, q_0, F) = (Q_1 \cup \{mid\} \cup Q_2, \delta, q_{01}, F_2)$  be its translated  $\epsilon$ -NFA using Definition 5.1. Then for any string  $w$ ,

3.1) if  $(e^*)$  accepts  $w$ , then there must exists a number  $n$  such that  $w \in (L \wedge n)$ . Now, lets do induction on  $n$ . **Base case:** when  $n = 0$ ,  $L \wedge 0 = \dots$

3.2) if  $M$  accepts  $w$ , ...

3.3) combining 3.1 and 3.2, we have  $L(e_1^*) = L(M)$ .  $\square$

## 4.6 Non-deterministic Finite Automata

**Definition 6.1** A NFA is a 5-tuple  $M = (Q, \Sigma, \delta, q_0, F)$ , where

1.  $Q$  is a finite set of states;
2.  $\Sigma$  is the set of alphabets;
3.  $\delta$  is a mapping from  $Q \times \Sigma$  to  $\mathcal{P}(Q)$  which defines the behaviour of the automata;
4.  $q_0$  in  $Q$  is the initial state;
5.  $F \subseteq Q$  is the set of accepting states.

Listing 5: NFA

```
record NFA : Set1 where
  field
    Q      : Set
    δ      : Q → Σ → DecSubset Q
    q0    : Q
    F      : DecSubset Q
    Q?     : DecEq Q
    |Q|-1  : ℕ
    It     : Vec Q (suc |Q|-1)
    ∀q∈It  : (q : Q) → (q ∈V It)
    unique : Unique It
```

The set of alphabets  $\Sigma : Set$  is passed to the file parameters. Together with  $Q$ ,  $\delta$ ,  $q_0$  and  $F$ , these five fields correspond to the 5-tuple  $\epsilon$ -NFA.  $Q?$  is the decidable equality of  $Q$ .  $|Q| - 1$  is the number of states - 1. ' $It$ ' is a vector of length  $|Q|$  containing all the states in  $Q$ .  $\forall q \in It$  is a proof that all states in  $Q$  are also in the vector ' $It$ '.  $unique$  is a proof that there is no repeating elements in ' $It$ '. These extra fields are important when computing  $\epsilon$ -closures, we will look into them again later in more details.

Now, we want to define the set of strings  $\Sigma^*$  accepted by a given NFA. However, before we can do this, we have to define some operations.

**Definition 6.2** A configuration is a pair  $Q \times \Sigma^*$ .

**Definition 6.3** A move by an  $\epsilon$ -NFA  $N$  is represented by a binary function  $\vdash$  on configurations. We say that  $(q, aw) \vdash (q', w)$  for all  $w$  in  $\Sigma^*$  if and only if  $q' \in \delta(q, a)$  where  $a \in \Sigma$ .

$$\begin{aligned} \_ \vdash \_ & : (\mathbb{Q} \times \Sigma \times \Sigma^*) \rightarrow (\mathbb{Q} \times \Sigma^*) \rightarrow \text{Set} \\ (q, a, w) \vdash (q', w') = w & \equiv w' \times q' \in^d \delta \ q \ a \end{aligned}$$

**Definition 6.4** We say that  $C \vdash^0 C'$  if and only if  $C = C'$ . We say that  $C_0 \vdash^k C_k$  for any  $k \geq 1$  if and only if there exists a chain of configurations  $C_1, C_2, \dots, C_{k-1}$  such that  $C_i \vdash C_{i+1}$  for all  $0 \leq i \leq k$ .

$$\begin{aligned} \_ \vdash^k \_ & : (\mathbb{Q} \times \Sigma^*) \rightarrow \mathbb{N} \rightarrow (\mathbb{Q} \times \Sigma^*) \rightarrow \text{Set} \\ (q, w) \vdash^k \text{zero} - (q', w') & \\ & = q \equiv q' \times w \equiv w' \\ (q, w) \vdash^k \text{suc } n - (q', w') & \\ & = \exists [ p \in \mathbb{Q} ] \exists [ a \in \Sigma ] \exists [ u \in \Sigma^* ] \\ & (w \equiv a :: u \times (q, a, u) \vdash (p, u) \times (p, u) \vdash^k n - (q', w')) \end{aligned}$$

**Definition 6.5** We say that  $C \vdash^* C'$  if and only if there exists a number of chains  $n$  such that  $C \vdash^n C'$ .

$$\begin{aligned} \_ \vdash^* \_ & : (\mathbb{Q} \times \Sigma^*) \rightarrow (\mathbb{Q} \times \Sigma^*) \rightarrow \text{Set} \\ (q, w) \vdash^* (q', w') & = \exists [ n \in \mathbb{N} ] (q, w) \vdash^k n - (q', w') \end{aligned}$$

**Definition 6.6** For any string  $w$ , it is accepted by an NFA  $N$  if and only if there exists a chain of configurations from  $q_0, w$  to  $q, \epsilon$  where  $q \in F$ .

**Definition 6.7** The language accepted by an NFA is given by the set  $\{ w \mid \exists q \in F. (q_0, w) \vdash^* (q, \epsilon) \}$ .

$$\begin{aligned} L^N & : \text{NFA} \rightarrow \text{Language} \\ L^N \text{ nfa} & = \lambda w \rightarrow \exists [ q \in \mathbb{Q} ] (q \in^d F \times (q_0, w) \vdash^* (q, [])) \end{aligned}$$

## 4.7 Removing $\epsilon$ -transitions

...

## 4.8 Deterministic Finite Automata

...

## 4.9 Powerset Construction

...

#### 4.10 Minimal DFA

...

#### 4.11 Minimising DFA

...

## 5 Further Extensions

Myhill-Nerode Theorem, Pumping Lemma

## 6 Evaluation

### 6.1 Correctness and Readability

According to Geuvers [8], a proof has two major roles: 1) to convince the readers that a statement is correct and 2) to explain why a statement is correct. The first role requires the proof to be convincing. The second role requires the proof to be able to give an intuition of why the statement is correct. In the paragraphs below, we will discuss these two criteria.

**Correctness** Traditionally, when a mathematician submits the proof of his/her concepts, a group of other mathematicians will evaluate and check the correctness of the proof. Alternatively, if we formalise the proof in a proof assistant, the proof will be checked automatically by the compiler. The only difference is that we are now relying on the compiler and the machine that runs the compile rather than the group of mathematicians. Therefore, if the compiler and the machine works properly, then any formalised proof that can be compiled without errors are said to be correct. In our case, we have the type checker and the termination checker in Agda to serve the purpose. Furthermore, a proof is consist of smaller reasoning steps. We can say that the a proof is correct if and only if all the reasoning steps within the proof are correct. When writing proofs in paper, we usually omit the proofs of some obvious lemmas and this sometime leads to mistakes. However, in Agda, we have to provide the proofs of every lemma we have used explicitly. Therefore, the correctness of an Agda proof always depend on the correctness of the smaller reasoning steps that it contains.

**Readability** The second purpose of a proof is to explain why a certain statement is correct. Let us consider the following code snippet extracted from **Correctness/RegExp-eNFA.agda**.

```
lem3 : ∀ we n q1
  → (q0 , we) ⊢k suc n = (inj q1 , [])
  → Σ[ n1 ∈ N ] Σ[ ue ∈ Σe* ] (toΣ* we = toΣ* ue × (inj q01 , ue) ⊢k n1 = (inj q1 , []))
lem3 _ zero q1 (inj .q1 , α _ , .[] , refl , (refl , ()), (refl , refl))
lem3 _ zero q1 (inj .q1 , E , .[] , refl , (refl , prf1) , (refl , refl)) with Q1? q1 q01
lem3 _ zero .q01 (inj .q01 , E , .[] , refl , (refl , prf1) , (refl , refl)) | yes refl = zero , [] , (refl , (refl , refl))
lem3 _ zero q1 (inj .q1 , E , .[] , refl , (refl , ()), (refl , refl)) | no p=q01
lem3 _ (suc n) q1 (init , α _ , ue , refl , (refl , ()), prf2)
lem3 _ (suc n) q1 (init , E , ue , refl , (refl , prf1) , prf2) = lem3 ue n q1 prf2
lem3 _ (suc n) q1 (inj p , α _ , ue , refl , (refl , ()), prf2)
lem3 _ (suc n) q1 (inj p , E , ue , refl , (refl , prf1) , prf2) with Q1? p q01
lem3 _ (suc n) q1 (inj .q01 , E , ue , refl , (refl , prf1) , prf2) | yes refl = suc n , ue , refl , prf2
lem3 _ (suc n) q1 (inj p , E , ue , refl , (refl , ()), prf2) | no p=q01
```

The above code is a proof that if  $w^e$  can take  $q_0$  to another state  $inj\ q_1$  in an  $\epsilon$ -NFA translated from a regular expression  $e^*$ , then there exists a number  $n_1$  and an  $\epsilon$ -string  $u^e$  that will take  $inj\ q_{01}$  to  $inj\ q_1$  where  $q_{01}$  is the start state of  $e$  and  $q_1$  is a state in  $e$ . There are several techniques used in the proof including induction on natural numbers and case analysis of state



comparison. However, by just looking at the function body, we can hardly understand the proving process. In conclusion, a computer proof is very inadequate on this purpose and thus we still need to outline the concept of the proof using natural language.

## 6.2 Different choices of representations

When writing computer proofs, we are usually required to provide a concrete representation for abstract mathematical objects. The consequence is that different representations will lead to different formalisations of the theorems and thus contribute to the easiness or difficulty in completing the proofs. In our project, we have made several decisions over the representations of different objects. In the following parts, we will discuss their consequences.

**The set of states ( $Q$ ) and its subsets** As we mentioned in section 3, Firsov and Uustalu [7] represented the set of states  $Q$  and its subsets as column matrices. However, this definition looks unnatural compare to the actual mathematical definition. Therefore, at the beginning of our project, we tried to avoid the vector representation. In our approach, the set of states are represented as a data type in Agda, i.e.  $Q : Set$ , and its subsets are represented as unary functions on  $Q$ , i.e.  $DecSubset A = A \rightarrow Bool$ .

Our definition allows us to finish the proofs in **Correctness/RegExp-eNFA.agda** without having to manipulate matrices. The proofs look much more natural compare to that in [7]. However, the problem arises when we need to iterate the set of states and its subsets when we are computing the  $\epsilon$ -closure. However, by using this definition, there is no possible way for us to iterate the sets. Therefore, we still need to include ***It*** – a vector containing all the states in  $Q$ , in the definition of automata. By using ***It***, we can iterate the subsets of  $Q$  by applying the function  $Q \rightarrow Bool$  to all the elements in ***It***. Note that the vector ***It*** is actually the vector representation of the set of states and therefore, we cannot avoid it.

**The language accepted by regular expression** At first, we defined the language accepted by regular expression as a decidable subset, i.e.  $L^R : RegExp \rightarrow (\Sigma^* \rightarrow Bool)$ .

## 6.3 Problems arise

Let us recall the definition of reachable states in **Translation/DFA-MDFA.agda**.

```
-- Reachable from  $q_0$ 
Reachable :  $Q \rightarrow Set$ 
Reachable  $q = \Sigma[ w \in \Sigma^* ] (q_0, w) \vdash^* (q, [])$ 

data  $Q^R : Set$  where
  reach :  $\forall q \rightarrow Reachable\ q \rightarrow Q^R$ 
```

We say that a state  $q$  is reachable if and only if there exists a string  $w$  that can take  $q_0$  to  $q$ . Therefore, the set  $Q^R$  should contain all and only the reachable states in  $Q$ . However, there may exist more than one reachability proof for a single state. This implies that there may be more than one element in  $Q^R$  having the same state. Therefore,  $Q^R$  may be larger than the original set  $Q$  or even worse, it may be infinite. This leads to a problem when we construct a new DFA using  $Q^R$  as the set of states. If  $Q^R$  is infinite, then we have no possible way to iterate the set in quotient construction. Even if the set  $Q^R$  is finite, the computational cost will be much higher. This is also one of the reasons why we cannot finish the quotient construction. However, surprisingly, this has no effects to the proof of  $L(DFA) = L(MDFA)$  because we can provide an equality relation of states in the record **DFA**. In the translation from DFA to MDFA, we defined the equality relation as follow: any two states in  $Q^R$  are equal if and only if the input states are equal. Therefore, two elements with same state but different reachability proofs are considered to be the same state in the new DFA.

One possible way to solve the problem is to re-define the reachability of a state such that any reachable state will have a unique reachability proof. For example, we can sort the reachability proofs according to alphabetical order of the string  $w$  and choose the proof with shortest string as the representative. This also requires a proof that the chosen proof is unique. Another solution is to use homotopy type to declare the set  $Q^R$ . This type allows us to group different reachability proofs into a single element such that every state will only appear once in  $Q^R$ .

## 7 Conclusion

Recap what is done

Is it feasible to write computer proof in practice?

## Bibliography

- [1] Alexandre Agular and Bassel Manna. Regular expressions in agda, 2009.
- [2] Alfred V. Aho and Jeffery V. Ullman. The theory of parsing, translation and compiling. volume i: Parsing, 1972.
- [3] Jeremy Avigad. Classical and constructive logic, 2000.
- [4] Ana Bove and Peter Dybjer. Dependent types at work. In *LERNET 2008. LNCS*, pages 57–99. Springer, 2009.
- [5] Haskell Curry. Functionality in combinatory logic, 1934.
- [6] Nicolaas de Bruijn. Automath, a language for mathematics, 1968.
- [7] Denis Firsov and Tarmo Uustalu. Certified parsing of regular languages, 2013.
- [8] Herman Geuvers. Proof assistants: history, ideas and future, 2009.
- [9] William A. Howard. The formulæ-as-types notion of construction. In Philippe De Groote, editor, *The Curry-Howard Isomorphism*. Academia, 1995.
- [10] Ivor. <https://eb.host.cs.st-andrews.ac.uk/ivor.php>. Accessed: 28th March 2016.
- [11] Per Martin-Löf. Intuitionistic type theory, 1984.
- [12] Bengt Nordström, Kent Peterson, and Jan M. Smith. *Programming in Martin-Löf’s Type Theory*. Oxford Clarendon Press, 1990.
- [13] Ulf Norell. Towards a practical programming language based on dependent type theory, 2007.
- [14] Ulf Norell and James Chapman. Dependently typed programming in agda. In *Lecture Notes from the Summer School in Advanced Functional Programming*, 2008.
- [15] The agda wiki. <http://wiki.portal.chalmers.se/agda/pmwiki.php?n=Main.HomePage>. Accessed: 12th March 2016.
- [16] Simon Thompson. *Type Theory and Functional Programming*. Addison Wesley Longman Publishing Co., Inc., Redwood City, CA, USA, 1991.
- [17] Philip Wadler. Propositions as types. *Commun. ACM*, 58(12):75–84, November 2015.

## Appendices

Agda Code?