#### Constructive Formalization of Regular Languages

Jan-Oliver Kaiser Advisors: Christian Doczkal, Gert Smolka Supervisor: Gert Smolka

UdS

November 7, 2012

Jan-Oliver Kaiser (UdS)

## Table of Contents

- 1 Motivation
- 2 Recap

Regular Expressions Finite Automata Finite Automata Regular Expression to Finite Automaton

- 3 Previous Talk Finite Automaton to Regular Expression
- 4 Myhill-Nerode Theorem Formalization of Equivalence Classes Myhill Relations and Nerode Relations Minimizing Equivalence Classes
- 5 Summary

Jan-Oliver Kaiser (UdS)

#### Motivation

- **Goal:** Build an extensive, elegant, constructive formalization of regular languages that includes
  - 1 regular expressions,
  - 2 the decidability of equivalence of regular expressions,
  - 3 finite automata,
  - 4 and the Myhill-Nerode theorem.
- How:  $\operatorname{Coq}$  with  $\operatorname{SSReflect}$ . No sacrifices for performance.

## **Regular Expressions**

Definition

• We use extended Regular Expressions (RE) over an alphabet  $\Sigma$ :

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

$$\begin{split} \mathcal{L}(\emptyset) &= \emptyset & \mathcal{L}(r^*) = \mathcal{L}(r)^* \\ \mathcal{L}(\varepsilon) &= \{\varepsilon\} & \mathcal{L}(r+s) = \mathcal{L}(r) \cup \mathcal{L}(s) \\ \mathcal{L}(.) &= \Sigma & \mathcal{L}(r\&s) = \mathcal{L}(r) \cap \mathcal{L}(s) \\ \mathcal{L}(a) &= \{a\} & \mathcal{L}(rs) = \mathcal{L}(r) \cdot \mathcal{L}(s) \\ \mathcal{L}(\neg r) &= \Sigma^* \backslash \mathcal{L}(r) \end{split}$$

- Implementation taken from Coquand and Siles<sup>1</sup>. This saved us quite a bit of work.
- $\approx$  200 lines of code including an implementation of regular languages and lots of useful lemmas.

Jan-Oliver Kaiser (UdS)

<sup>&</sup>lt;sup>1</sup>Thierry Coquand and Vincent Siles. "A Decision Procedure for Regular Expression Equivalence in Type Theory". In: *CPP*. 2011, pp. 119–134.

## Finite Automata

Definition

- Our Finite Automata (FA) over an alphabet  $\boldsymbol{\Sigma}$  consist of
  - 1 a set of states Q,
  - 2 a starting state  $s \in Q$ ,
  - **3** a set of final states  $F \subseteq Q$ ,
  - 4 a transition relation  $\delta \in Q \times \Sigma \times Q$ .
- Two types: one for non-deterministic FA ( $\delta$  may be non-functional), one for deterministic FA ( $\delta$  is functional).
- For our deterministic FA,  $\delta$  is also total and, thus, a function. This helped, but maybe not a lot.

Jan-Oliver Kaiser (UdS)

#### Finite Automata

Formalization

- Our definition is very close to the textbook definition.
- pprox 120 lines of code including some general lemmas for later proofs.

## $\mathsf{RE} \implies \mathsf{FA}$

- Structure of proof given by inductive definition of RE.
- Plan: Construct FA for every RE constructor.
- Sounds simple enough ..
- ..  $\approx 800$  lines of code.
- $\approx 100$  lines of that needed for equivalence of DFA and NFA.
- Another  $\approx 100$  lines of code for extended regular expressions.
- This is a candidate for improvement.

## $FA \implies RE$

- We use the "Transitive Closure method", Kleene's original proof<sup>2</sup>.
- This method recursively builds a regular expression  $R_{x,y}^X$  that recognizes words whose runs starting in x only pass through states in X and end in y.
- The previous version constructed R<sup>k</sup><sub>x,y</sub> which translates to R<sup>{z|#(z)<k}</sup><sub>x,y</sub> where # is an ordering on Q.
- Instead of **nat**, we now recurse on the size of a **finite subset** of  $Q^{3}$ .
- This also avoids cumbersome conversions from nat to SSREFLECT's ordinals and, finally, to states.

<sup>&</sup>lt;sup>2</sup>S. C. Kleene. "Representation of Events in Nerve Nets and Finite Automata". In: Automata Studies (1965).

<sup>&</sup>lt;sup>3</sup>Dexter Kozen. Automata and computability. Undergraduate texts in computer science. Springer, 1997, pp. I–XIII, 1–400. ISBN: 978-0-387-94907-9.

 $FA \implies RF$ 

Proof Outline

- We introduce a decidable language L<sup>X</sup><sub>x,y</sub> that directly encodes the idea behind R<sup>X</sup><sub>x,y</sub> as a boolean predicate.
- The proof of FA  $\implies$  RE consists of three steps:
  - **1** We show that  $L_{x,y}^{X}$  respects the defining recursive equation of  $R_{x,y}^{X}$ .
  - 2 We show that for all  $w \in \Sigma^*$ ,  $w \in L^X_{x,y} \iff w \in R^X_{x,y}$ .
  - 3 We show that  $\bigcup_{f \in F} L_{s,f}^Q = \mathcal{L}(A)$  and thus  $\sum_{f \in F} R_{s,f}^Q = \mathcal{L}(A)$ .

## $\mathsf{FA} \implies \mathsf{RE}$

- After some restructuring:  $\approx 550$  lines of code,  $\approx 150$  of which are general infrastructure.
- Previous version:  $\approx$  800 lines of code, much harder to read.

```
Lemma L_split k' i j a w:

let k := k_ord k' in

(a::w) \in L^k'.+1 i j ->

(a::w) \in L^k' i j \/

exists w1, exists w2,

a::w = w1 ++ w2 /\

w1 != [::] /\

w1 \in L^k' i (enum_val k) /\

w2 \in L^k'.+1 (enum_val k) j.
```

Figure: Previous and current version of the same lemma

Jan-Oliver Kaiser (UdS)

- It turns out that there are two different concepts: Myhill relations and the Nerode relation.
- We also consider a related characterization: weak Nerode relations.
- All these are equivalence relations which we require to be of finite index, i.e. to have a finite number of equivalence classes.
- $\bullet\,$  However,  $\mathrm{Coq}$  has no notion of quotient types.

Equivalence Relations of Finite Index

- We use functions of finite domain to represent equivalence relations of finite index.
- Think of the domain as the set of equivalence classes.
- We also need to have a representative of every equivalence class. Thus, we require surjectivity.

```
Record Fin_Eq_Cls :=
    { fin_type : finType;
        fin_func :> word -> fin_type;
        fin_surjective : surjective fin_func }.
```

• With constructive choice, we can then give a canonical representative of every every equivalence class.

```
Definition cr (f: Fin_Eq_Cls) x := xchoose (fin_surjective f x).
```

Jan-Oliver Kaiser (UdS)

Myhill relations, weak Nerode relations, Nerode relation

- An equivalence relation  $\equiv$  of finite index is a  $Myhill\ relation^4$  for L iff
  - **1** = is right-congruent, i.e.  $\forall u, v \in \Sigma^*$ .  $\forall a \in \Sigma$ .  $u \equiv v \implies ua \equiv va$ ,
  - 2 and  $\equiv$  refines L, i.e.  $\forall u, v \in \Sigma^*$ .  $u \equiv v \implies (u \in L \iff v \in L)$ .
- An equivalence relation  $\equiv$  of finite index is a weak Nerode relation for L iff

$$\forall u, v \in \Sigma^*. \ u \equiv v \implies \forall w \in \Sigma^*. \ (uw \in L \iff vw \in L).$$

• The **Nerode relation**<sup>5</sup>  $\doteq_L$  for L is defined such that

$$\forall u, v \in \Sigma^*. \ u \doteq_L v \iff \forall w \in \Sigma^*. \ (uw \in L \iff vw \in L).$$

Jan-Oliver Kaiser (UdS)

<sup>&</sup>lt;sup>4</sup> John R. Myhill. *Finite Automata and the Representation of Events*. Tech. rep. WADC TR-57-624. Wright-Paterson Air Force Base, 1957.

<sup>&</sup>lt;sup>5</sup>Anil Nerode. "Linear Automaton Transformations". In: *Proceedings of the American Mathematical Society* 9.4 (1958), pp. 541–544.

Formalization of Myhill, weak Nerode and Nerode relation

- For all three relations, we build a record that consists of an equivalence relation of finite index and proofs of the properties of the respective relation.
- Example:

```
Record Myhill_Rel :=
```

{ myhill\_func :> Fin\_Eq\_Cls; myhill\_congruent : right\_congruent myhill\_func ; myhill\_refines : refines myhill\_func }.

- Our version of the Myhill-Nerode theorem states that the following are equivalent
  - 1 language L is accepted by a DFA,
  - 2 we can construct a Myhill relation for L,
  - 3 we can construct a weak Nerode relation for L,
  - 4 the Nerode relation for L is of finite index.
- (1)  $\implies$  (2) is easy. (Map word w to the last state of its run on the automaton.)
- (2)  $\implies$  (3) is a trivial inductive proof.
- (4)  $\implies$  (1) is also straight-forward. (Use equivalence classes as states.)

Jan-Oliver Kaiser (UdS)

 $(3)\implies (4)$ 

- A weak Nerode relation (given as a function f) has at least as many equivalence classes as the Nerode relation.
- Some of them are redundant w.r.t. to L, i.e. we may have equivalence classes s.t.

 $\forall uv \in \Sigma^*$ .  $fu \neq fu \land \forall w \in \Sigma^*$ .  $(uw \in L \iff vw \in L)$ .

- Our goal is to merge these equivalence classes.
- In fact, we construct an equivalence relation on these equivalence classes.
- Equivalence classes are contained in that relation iff they are equivalent w.r.t. to L.

Jan-Oliver Kaiser (UdS)

Finding equivalent equivalence classes

- We use the table-filling algorithm<sup>6</sup> for minimizing DFA.
- Our fixed-point algorithm finds all equivalence classes whose representatives are distinguishable w.r.t. to L.
- The remaining equivalence classes are then equivalent w.r.t. to L.
- Start:  $\{(x,y) \mid cr(x) \notin L \iff cr(y) \in L\}$ .
- Step: if the previous result is d, the new result is  $d \cup \{(x, y) \mid \exists a. (f(cr(x)a), f(cr(y)a)) \in d\}.$
- Due to finiteness of the domain and monotonicity of the algorithm, it has a fixed point which we can compute.

<sup>&</sup>lt;sup>6</sup>D.A. Huffman. "The synthesis of sequential switching circuits". In: *Journal of the Franklin Institute* 257.3 (1954), pp. 161 –190.

Proof Outline

- We construct a function *f<sub>min</sub>* that maps every equivalence class to the set of equivalence classes equivalent to it w.r.t. L.
- We then show that  $f_{min}$  is surjective and encodes an equivalence relation of finite index.
- Finally, we show that  $f_{min}$  is equivalent to the Nerode relation.

- The lemmas of this chapter are rather abstract, which makes for nice and short statements.
- The proofs also received more refinement than the other chapters. They are very concise.
- The entire chapter consists of  $\approx$  430 lines of code.

## Summary

- All in all we have  $\approx 2100$  lines of code.
- Mostly very close to the mathematical definitions.
- Code produced in the beginning of the project might be improved quite a bit.

# Thank you for your attention!

Jan-Oliver Kaiser (UdS)

Constr. Formalization of Reg. Languages

November 7, 2012 21 / 21