# Constructive Formalization of Regular Languages

## Jan-Oliver Kaiser

Advisors: Christian Doczkal, Gert Smolka
Supervisor: Gert Smolka

# Contents

1. Motivation
2. Quick Recap
3. Previous work
4. Our development
5. Roadmap

# Motivation

We want to develop an elegant formalization of regular languages in Coq based on finite automata.

There are several reasons for choosing this topic and our specific approach:

- Strong interest in formalizations in this area.

- Few formalizations of regular languages in Coq, most of them very long or incomplete.

- Most formalizations avoid finite automata in favor of regular expressions. Regular expressions (with Brzozowski derivatives) lead to more complex but also more performant algorithms.

# Quick Recap

We use extended **regular expressions** (regexp):

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

- $\mathcal{L}(\emptyset) := \{\}$
- $\mathcal{L}(\varepsilon) := \{\varepsilon\}$
- $\mathcal{L}(a) := \{a\}$
- $\mathcal{L}(rs) := \mathcal{L}(r) \cdot \mathcal{L}(s)$
- $\mathcal{L}(r+s) := \mathcal{L}(r) \cup \mathcal{L}(s)$
- $\mathcal{L}(r\&s) := \mathcal{L}(r) \cap \mathcal{L}(s)$
- $\mathcal{L}(r^{*}) := \mathcal{L}(r)^{*}$
- $\mathcal{L}(\neg r) := \overline{\mathcal{L}(r)}$

**Quick Recap**

## **Derivatives of Regular Expressions** (1964), *Janusz Brzozowski*:

- der a $\emptyset = \emptyset$

- der a $\varepsilon = \emptyset$

- der a b = if a = b then $\varepsilon$ else $\emptyset$

- der a (r s) = if $\delta(r)$ then (der a s) + ((der a r) s) else (der a r) s
  with $\delta(r) = true \Leftrightarrow \varepsilon \in \mathcal{L}(r)$.

- der a (r + s) = (der a r) + (der a s)

- der a (r & s) = (der a r) & (der a s)

- der a (r*) = (der a r) r*

- der a ( ¬r) = ¬(der a r)

**Theorem**: $w \in \mathcal{L}(r)$ if and only if the derivative of r with respect to $w_1 .. w_{|w|}$ accepts $\varepsilon$.

**Quick Recap**

Regular languages are also exactly those languages accepted by **finite automata** (FA).

Our definition of FA over an alphabet $\Sigma$ :

- The finite set of states Q
- The initial state $s_0 \in$ Q
- The (decidable) transition relation $\Delta \in$ (Q, $\Sigma$, Q)

  Deterministic FA: $\Delta$ is functional and **total**.
- The set of finite states F, F $\sqsubseteq$ Q

Let A be a FA.
$$\mathcal{L}(A) := \left\{ w \mid \exists s_1, \ldots s_{|w|} \in Q \ s.t. \ \forall i : 0 < i \le n \rightarrow (s_{i-1}, w_i, s_i) \in \Delta \right\}$$

Finally, regular languages are also characterized by the Myhill-Nerode theorem.

First, we define an equivalence relation on L:

$$x \; R_L \; y \; := \; \forall z, \; x{\cdot}z \; \in \; L \; \Leftrightarrow \; y{\cdot}z \; \in \; L$$

**Myhill-Nerode theorem**: L is regular if and only if $R_L$ divides L into a finite number of equivalence classes.

**Quick Recap**

# Previous work

- Constructively formalizing automata theory (2000)

  *Robert L. Constable, Paul B. Jackson, Pavel Naumov, Juan C. Uribe*

  **PA**: Nuprl

  The first constructive formalization of MH.

  Based on **FA**.

- Proof Pearl: Regular Expression Equivalence and Relation Algebra (2011)

  *Alexander Krauss, Tobias Nipkow*

  **PA**: Isabelle

  Based on **derivatives of regexps**. No proof of termination.

- Deciding Kleene Algebras in Coq (2011)

  *Thomas Braibant, Damien Pous*

  **PA**: Coq

  Based on **FA**, matrices. Focus on performance.

- A Decision Procedure for Regular Expression Equivalence in Type Theory (2011)

  *Thierry Coquand, Vincent Siles*

  **PA**: Coq

  Based on **derivatives of regexps**.

**Previous work**

- A Formalisation of the Myhill-Nerode Theorem based on Regular Expressions (Proof Pearl) (2011)

  *Chunhan Wu, Xingyuan Zhang, Christian Urban*

  **PA**: Isabelle

  The first proof of MH based on **derivatives of regexps**.

- Deciding Regular Expressions (In-)Equivalence in Coq (2011)

  Nelma Moreira, David Pereira, Simão Melo de Sousa

  **PA**: Coq

  Based on Krauss, Nipkow. Proof of termination.

**Previous work**

## Our Development

- We want to focus on elegance, not performance.
- Our main goals are MH and the decidability of regexp equivalence.
- We use finite automata.

  They are not at all impractical. (Partly thanks to Ssreflect's finType)

## Ssreflect

- Excellent support for all things boolean.
- Finite types with all necessary operations and closure properties. (very useful for alphabets, FA states, etc.)
- Lots and lots of useful lemmas and functions.

## Finite automata

DFA and NFA without e-transitions.

- DFA to prove closure under ∪, ∩, and ¬.
- NFA to prove closure under · and ∗.

Also proven: NFA ⇔ DFA.

This gives us: regexp ⇒ FA.

# Roadmap

1. Emptiness test on FA ( $\emptyset(A) := \mathcal{L}(A) = \emptyset$ )

2. FA $\Rightarrow$ regexp

3. Dedicedability of regexp equivalence using regexp $\Rightarrow$ FA, (2) and (1):

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

   $$\Leftrightarrow$$

   $$\emptyset(\mathcal{A}(r) \cap \overline{\mathcal{A}(s)}) \wedge \emptyset(\overline{\mathcal{A}(r)} \cap \mathcal{A}(s))$$

4. Finally, we want to prove the MH theorem

With this we'll have an extensive formalization of regular languages including regular expressions, FA and MH and all corresponding equivalences.

**Roadmap**