Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

# Constructive Formalization of Regular Languages

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

April 27, 2012

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

1 Goal

2 State Removal

3 Brzozowski Algebraic Method

4 Transitive Closure

5 Our Approach

**Goal:**

Find simple proofs for the decidability of regular expression equivalence and the Myhill-Nerode theorem.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Current Goal
Roadmap

Roadmap:

1. RE $\Rightarrow$ FA (**DONE**)

2. Emptiness test on FA (**Easy**)

3. RE equivalence (**Follows from 1 and 2**)

4. FA $\Rightarrow$ RE (**Work in progress**)

5. Myhill-Nerode

Goal
**Recap**
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

**Definitions:**

- We use extended regular Expressions (RE):

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

$$\textbf{Def.:} \quad \underset{x \in X}{+}\, r_x := r_{x_0} + ... + r_{x_{|X|-1}}$$

- Our Finite Automata (FA) are

$$(\Sigma, Q, q_0, F, \delta)$$

(The transition relation $\delta$ of deterministic FA is total).

## Finite automata to regular expressions

- Converting REs to FAs is straight-forward and there is really only one algorithm (with slight variations): We mirror the constructors of RE in operations on FAs.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Difficulty
Overview

**Finite automata to regular expressions**

- Converting REs to FAs is straight-forward and there is really only one algorithm (with slight variations):
  We mirror the constructors of RE in operations on FAs.

- Converting FAs to REs is complicated and there are at least three algorithms found in textbooks.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Difficulty
Overview

## Finite automata to regular expressions

- Converting REs to FAs is straight-forward and there is really only one algorithm (with slight variations): We mirror the constructors of RE in operations on FAs.

- Converting FAs to REs is complicated and there are at least three algorithms found in textbooks.

### Why is that?

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Difficulty
Overview

**My intuition:**

- Converting REs to FAs is done by structural recursion on a **tree**. The result is a **flat structure**.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Difficulty
Overview

**My intuition:**

- Converting REs to FAs is done by structural recursion on a **tree**. The result is a **flat structure**.

- Converting FAs to REs **can not be done** by structural recursion. There is **no recursive structure** in FAs.
  But somehow we need to construct a **tree** of REs.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Difficulty
Overview

**Three methods (+ variations):**

1. Transitive Closure

2. State Removal
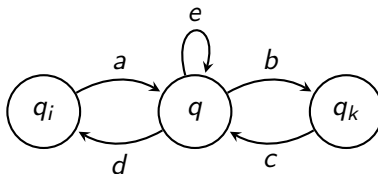
3. Brzozowski Algebraic Method

**State Removal**

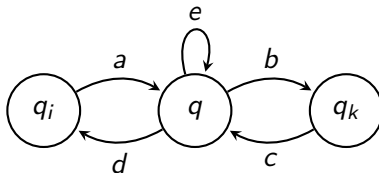**Given**: NFA $A = (\Sigma, Q, q_0, F, \delta)$.

**New concept**: Automata that have transitions labeled by RE.

**Idea**: Remove states until there are two or less states remaining. Update the remaining states' transitions by incorporating the "lost" paths.
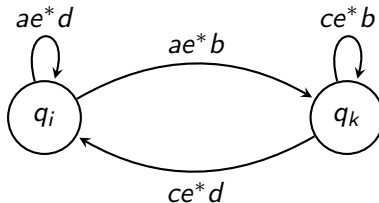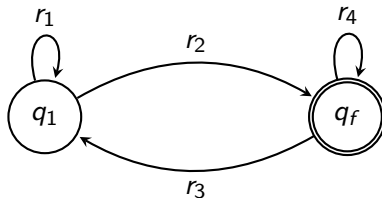
**Remove $q$ from**

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Approach
Algorithm
Caveats
Properties

**Remove $q$ from**



**to get**

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Approach
Algorithm
Caveats
Properties

**Repeat until $A$ is of this form:**



$$\Rightarrow \mathcal{L}(A) = \mathcal{L}(r_1^* r_2 (r_4 + r_3 r_1^* r_2)^*)$$
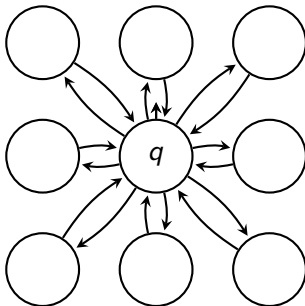
## Caveats

- It looks like we only need to update two edges.

#### Caveats

- It looks like we only need to update two edges.
  In reality, there can be $|Q| - 1$ states connected to $q$.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Approach
Algorithm
Caveats
Properties

**Caveats**

- It looks like we only need to update two edges.
  In reality, there can be $|Q| - 1$ states connected to $q$.

**Caveats**

- It looks like we only need to update two edges.
  In reality, there can be $|Q| - 1$ states connected to $q$.
- What about final states?

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Approach
Algorithm
Caveats
Properties

**Caveats**

- It looks like we only need to update two edges.
  In reality, there can be $|Q| - 1$ states connected to $q$.

- What about final states?

  1 Introduce a new final state without any outgoing edges.
  2 Introduce $\varepsilon$ transitions from all other final states to the final new state.
  3 Make all other states non-final.
  4 Never remove the new final state.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Approach
Algorithm
Caveats
Properties

**Formalization:**

- Requires a new kind of finite automaton that has RE transitions.
- Lots of details to consider.
- Induction on the number of states.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Approach
Algorithm
Properties

**Brzozowski Algebraic Method**

**Given**: FA $A = (\Sigma, Q, q_0, F, \delta)$.

**Idea**: Retrieve a RE for FA by solving a system of equations determined by $\delta$.

Goal
Recap
Finite automata to regular expressions    Approach
State Removal    **Algorithm**
Brzozowski Algebraic Method    Properties
Transitive Closure
Our Approach

**Construct system of equations**:

$$
\begin{aligned}
r_0 \quad &= \sum_{\substack{a \in \Sigma \\ 0 \le i < |Q|}} \{\, a\, r_i \mid (q_0, a, q_i) \in \delta \} \qquad (+\,\varepsilon \text{ if } r_0 \in F) \\
\vdots \quad &= \vdots \\
r_{|Q|-1} \quad &= \sum_{\substack{a \in \Sigma \\ 0 \le i < |Q|}} \{\, a\, r_i \mid (q_{|Q|-1}, a, q_i) \in \delta \} \quad (+\,\varepsilon \text{ if } r_{|Q|-1} \in F)
\end{aligned}
$$

$\Rightarrow \mathcal{L}(A) = \mathcal{L}(r_0)$

Goal
Recap
Finite automata to regular expressions          Approach
State Removal          **Algorithm**
**Brzozowski Algebraic Method**          Properties
Transitive Closure
Our Approach

Solve the system by substitution and **Arden's Lemma** which
states that for all regular languages X, Y and Z the equation

$$X = YX + Z \tag{1}$$

has the unique solution

$$X = Y^*Z \tag{2}$$

Goal
Recap
Finite automata to regular expressions    Approach
State Removal    Algorithm
Brzozowski Algebraic Method    Properties
Transitive Closure
Our Approach

**Formalization:**

- Requires a formalization of these equations and operations on them.
- We would need to prove Arden's Lemma.
- We would also need to prove that Arden's Lemma (and substitution) is enough to solve these systems of equations.

Goal
Recap
Finite automata to regular expressions    **Approach**
State Removal    Algorithm
Brzozowski Algebraic Method    Properties
**Transitive Closure**
Our Approach

**Transitive Closure**
**Given**: FA $A = (\Sigma, Q, q_0, F, \delta)$.
**Idea**: Construct regexps $r_f$ for every final states $f \in F$ s.t.
$r_f$ matches all words which A accepts with final state $f$.

$$\Rightarrow \mathcal{L}(A) = \mathcal{L}(\underset{f \in F}{+} r_f)$$

**How do we construct $r_f$?**

Goal
Recap
Finite automata to regular expressions    Approach
State Removal    **Algorithm**
Brzozowski Algebraic Method    Properties
**Transitive Closure**
Our Approach

We generalize the idea of $r_f$ to $R_{ij}^k$ which matches all words which lead from state $i$ to $j$ while passing only through states with index smaller than $k$.

1. Merge multiple edges between states to one unified edge.
2. Construct regexp $R_{ij}^k$ recursively:

$$R_{ij}^0 := \begin{cases} r & \text{if } i \neq j \land i \text{ has edge } r \text{ to j} \\ \varepsilon + r & \text{if } i = j \land i \text{ has edge } r \text{ to j} \\ \emptyset & \text{otherwise} \end{cases}$$

$$R_{ij}^k := R_{ik}^{k-1} R_{kk}^{k-1} R_{kj}^{k-1} + R_{ij}^{k-1}$$

$$\Rightarrow \mathcal{L}(A) = \mathcal{L}(\sum_{f \in F} r_f) = \mathcal{L}(\sum_{f \in F} R_{0f}^{|Q|})$$

Goal
Recap
Finite automata to regular expressions     Approach
State Removal     Algorithm
Brzozowski Algebraic Method     Properties
Transitive Closure
Our Approach

**Formalization:**

- Easier than the other methods.

- The recursive definition translates quite well.

- The details are quite challenging.

Goal
Recap
Finite automata to regular expressions     Approach
State Removal     Algorithm
Brzozowski Algebraic Method     Properties
Transitive Closure
Our Approach

**Formalization:**

- Easier than the other methods.

- The recursive definition translates quite well.

- The details are quite challenging.

    **This appears to be the simplest formalization.**

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Definitions
Lemmas

**Our Approach:**
There are different ways of formalizing $R_{ij}^k$ itself, especially its parameters. Most practical so far:
$k$ is of type nat, $i$ and $j$ are ordinals from $[0..|Q|-1]$.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Definitions
Lemmas

**Our Approach:**

There are different ways of formalizing $R_{ij}^k$ itself, especially its parameters. Most practical so far:

$k$ is of type nat, $i$ and $j$ are ordinals from $[0..|Q| - 1]$.

This gives us easy recursion and matching on $k$.

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Definitions
Lemmas

**Our Approach:**
There are different ways of formalizing $R_{ij}^k$ itself, especially its parameters. Most practical so far:
$k$ is of type nat, $i$ and $j$ are ordinals from $[0..|Q|-1]$.
This gives us easy recursion and matching on $k$.
**But** we have to use $min\, k\, (|Q|-1)$ to map k to the corresponding ordinal (and then to a state).

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Definitions
Lemmas

To get a FA counterpart to $R_{ij}^k$, we introduce $A_{ij}^k$ s.t.

$$\mathcal{L}(A_{ij}^k) = \mathcal{L}(R_{ij}^k).$$

$A_{ij}^k$ is similar to A. It has one additional state, which has all incoming edges of state $j$ and no outgoing edges. It only leaves states that are $i$ or $< k$. It only enters states less than $< k$ or the new state.
We can then show that

$$\mathcal{L}(\bigcup_{f \in F} A_{q_0 f}^{|Q|}) = \mathcal{L}(A).$$

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Definitions
Lemmas

**Lemmas:**

1. The language of an automaton is the union of the words with
   runs on $A$ that end in any final state, i.e.
   $\mathcal{L}(A) = \bigcup_{f \in F} \mathcal{L}(A_f)$, where $A_f$ accepts only in $f$. (Not
   difficult)

2. A path in $A_{ij}^{k+1}$ that is not in $A_{ij}^k$ can be decomposed into:
   - a sub path from index 0 to the first occurrence of k
   - a sub path from the first to the last occurrence of k
   - the remaining sub path from the last occurrence of k to the
     end.

   (Complex; requires a new operations on lists)

Goal
Recap
Finite automata to regular expressions
State Removal
Brzozowski Algebraic Method
Transitive Closure
Our Approach

Definitions
Lemmas

### More lemmas:

3. A path on $A$ which passes through a final state more than once is in $\mathcal{L}(A)^*$ (Relies on the same operations we need for lemma 2)

4. Any path on $A_{ij}^k$ that ends in the new accepting state can be mapped to an equivalent path that ends state j (and vice-versa).

5. Any path on $A_{ij}^k$ that does not end in the new accepting state is also a path on $A$.

# Thank you for your attention