

# POSIX Lexing with Derivatives of Regular Expressions (Proof Pearl)

Fahad Ausaf, Roy Dyckhoff and Christian Urban

King's College London, University of St Andrews

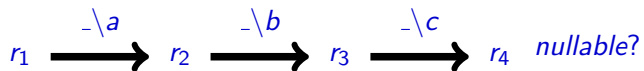
# Brzowski's Derivatives of Regular Expressions

Idea: If  $r$  matches the string  $c::s$ , what is a regular expression that matches just  $s$ ?

chars:	$0 \setminus c$	$\stackrel{\text{def}}{=} 0$
	$1 \setminus c$	$\stackrel{\text{def}}{=} 0$
	$d \setminus c$	$\stackrel{\text{def}}{=} \text{if } d = c \text{ then } 1 \text{ else } 0$
	$r_1 + r_2 \setminus c$	$\stackrel{\text{def}}{=} r_1 \setminus c + r_2 \setminus c$
	$r_1 \cdot r_2 \setminus c$	$\stackrel{\text{def}}{=} \text{if nullable } r_1$ $\text{then } r_1 \setminus c \cdot r_2 + r_2 \setminus c \text{ else } r_1 \setminus c \cdot r_2$
	$r^* \setminus c$	$\stackrel{\text{def}}{=} r \setminus c \cdot r^*$
strings:	$r \setminus []$	$\stackrel{\text{def}}{=} r$
	$r \setminus c::s$	$\stackrel{\text{def}}{=} (r \setminus c) \setminus s$

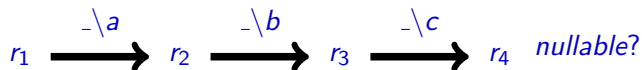
# Brzozowski's Matcher

Does  $r_1$  match string  $abc$ ?



# Brzozowski's Matcher

Does  $r_1$  match string  $abc$ ?

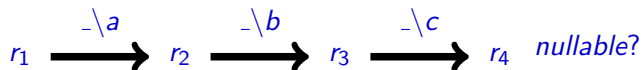


It leads to an elegant functional program:

$$\text{matches}(r, s) \stackrel{\text{def}}{=} \text{nullable}(r \backslash s)$$

# Brzozowski's Matcher

Does  $r_1$  match string  $abc$ ?



It leads to an elegant functional program:

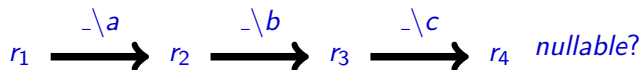
$$\text{matches}(r, s) \stackrel{\text{def}}{=} \text{nullable}(r \setminus s)$$

It is an easy exercise to formally prove (e.g. Coq, HOL, Isabelle):

$$\text{matches}(r, s) \text{ if and only if } s \in L(r)$$

# Brzozowski's Matcher

Does  $r_1$  match string  $abc$ ?



It leads to an elegant functional program:

$$\text{matches}(r, s) \stackrel{\text{def}}{=} \text{nullable}(r \setminus s)$$

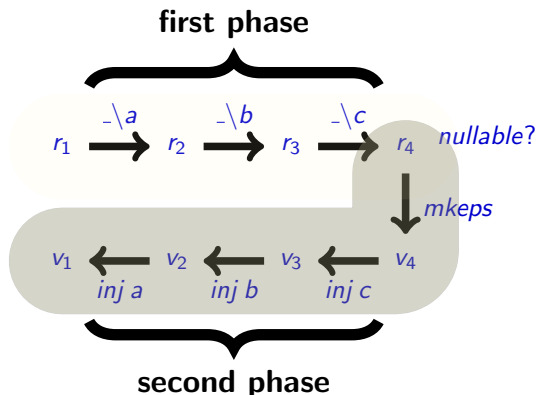
It is an easy exercise to formally prove (e.g. Coq, HOL, Isabelle):

$$\text{matches}(r, s) \text{ if and only if } s \in L(r)$$

**But Brzozowski's matcher gives only a yes/no-answer.**

# Sulzmann and Lu's Matcher

Sulzmann and Lu added a second phase in order to answer **how** the regular expression matched the string.



There are several possible answers for **how**: POSIX, GREEDY, ...

# Regular Expressions and Values

Regular expressions and their corresponding values (for how a regular expression matched a string):

$r$	::=	<b>0</b>	$v$	::=	<i>Empty</i>
		<b>1</b>			<i>Char(c)</i>
		$c$			<i>Seq(<math>v_1 \cdot v_2</math>)</i>
		$r_1 \cdot r_2$			<i>Left(v)</i>
		$r_1 + r_2$			<i>Right(v)</i>
		$r^*$			$[v_1, \dots, v_n]$



# POSIX Matching (needed for Lexing)

**Longest Match Rule:** The longest initial substring matched by any regular expression is taken as the next token.

**Rule Priority:** For a particular longest initial substring, the first regular expression that can match determines the token.

For example:  $r_{keywords} + r_{identifiers}$ :

- i f f o o \_ b l a
- i f \_ b l a

# Problems with POSIX

Grathwohl, Henglein and Rasmussen wrote:

*“The POSIX strategy is more complicated than the greedy because of the dependence on information about the length of matched strings in the various subexpressions.”*

Also Kuklewicz maintains a unit-test repository for POSIX matching, which indicates that most POSIX matchers are buggy.

[http://www.haskell.org/haskellwiki/Regex\\_Posix](http://www.haskell.org/haskellwiki/Regex_Posix)

## “Correctness” by Sulzmann and Lu

- Sulzmann & Lu’s idea is to order all possible answer such that they can prove the correct answer is the maximum
- The idea is taken from a GREEDY algorithm (and it works there)

# “Correctness” by Sulzmann and Lu

- Sulzmann & Lu’s idea is to order all possible answer such that they can prove the correct answer is the maximum
- The idea is taken from a GREEDY algorithm (and it works there)
- **But** we made no progress in formalising Sulzmann & Lu’s idea, because
  - transitivity, existence of maxima etc all fail to turn into real proofs
  - the reason: the ordering works only if ....
  - though we did find mistakes:

“How could I miss this? Well, I was rather careless when stating this Lemma :)

Great example how formal machine checked proofs (and proof assistants) can help to spot flawed reasoning steps.”

# “Correctness” by Sulzmann and Lu

- Sulzmann & Lu’s idea is to order all possible answer such that they can prove the correct answer is the maximum
- The idea is taken from a GREEDY algorithm (and it works there)
- **But** we made no progress in formalising Sulzmann & Lu’s idea, because
  - transitivity, existence of maxima etc all fail to turn into real proofs
  - the reason: the ordering works only if ....
  - though we did find mistakes:

“How could I miss this? Well, I was rather careless when stating this Lemma :)

Great example how formal machine checked proofs (and proof assistants) can help to spot flawed reasoning steps.”

“Well, I don’t think there’s any flaw. The issue is how to come up with a mechanical proof. In my world mathematical proof = mechanical proof doesn’t necessarily hold.”

# “Correctness” by Sulzmann and Lu

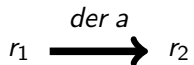
- Sulzmann & Lu’s idea is to order all possible answer such that they can prove the correct answer is the maximum
- The idea is taken from a GREEDY algorithm (and it works there)
- **But** we made no progress in formalising Sulzmann & Lu’s idea, because
  - transitivity, existence of maxima etc all fail to turn into real proofs
  - the reason: the ordering works only if ....
  - though we did find mistakes:

“How could I miss this? Well, I was rather careless when stating this Lemma :)

Great example how formal machine checked proofs (and proof assistants) can help to spot flawed reasoning steps.”

“Well, I don’t think there’s any flaw. The issue is how to come up with a mechanical proof. In my world mathematical proof = mechanical proof doesn’t necessarily hold.”

We want to match the string *abc* using  $r_1$

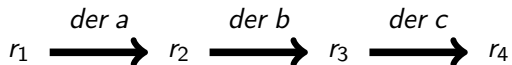


We want to match the string *abc* using  $r_1$

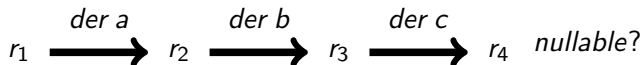




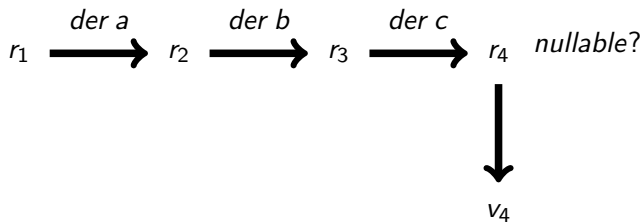
We want to match the string *abc* using  $r_1$



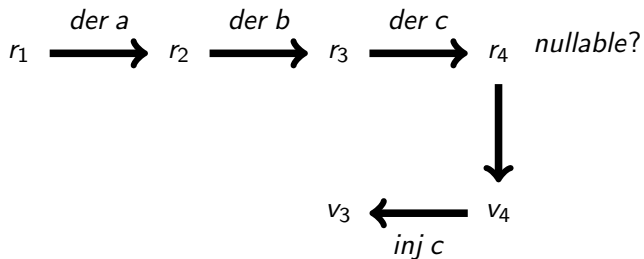
We want to match the string *abc* using  $r_1$



We want to match the string *abc* using  $r_1$

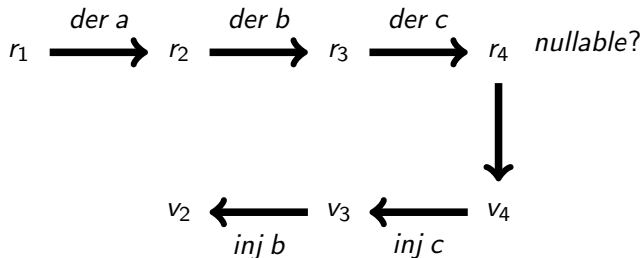


We want to match the string *abc* using  $r_1$



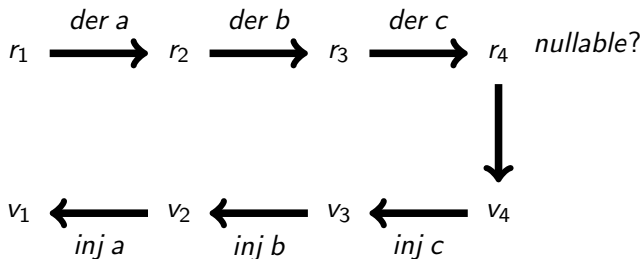
# Sulzmann and Lu Matcher

We want to match the string *abc* using  $r_1$

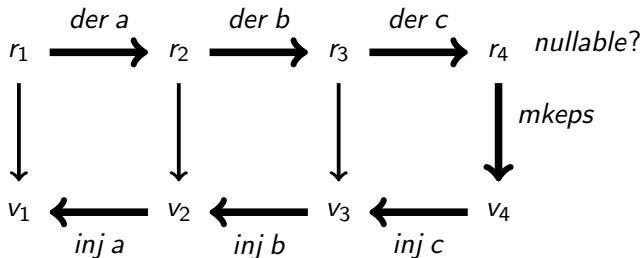


# Sulzmann and Lu Matcher

We want to match the string *abc* using  $r_1$



We want to match the string *abc* using  $r_1$



- Sulzmann: ... Let's assume  $v$  is not a *POSIX* value, then there must be another one ... contradiction.



- Sulzmann: ... Let's assume  $v$  is not a *POSIX* value, then there must be another one ... contradiction.
- Exists ?

$$L(r) \neq \emptyset \Rightarrow \exists v. \text{POSIX}(v, r)$$

- Sulzmann: ... Let's assume  $v$  is not a *POSIX* value, then there must be another one ... contradiction.
- Exists ?

$$L(r) \neq \emptyset \Rightarrow \exists v. \text{POSIX}(v, r)$$

- In the sequence case  $\text{Seq}(v_1, v_2) \succ_{r_1 \cdot r_2} \text{Seq}(v'_1, v'_2)$ , the induction hypotheses require  $|v_1| = |v'_1|$  and  $|v_2| = |v'_2|$ , but you only know

$$|v_1| @ |v_2| = |v'_1| @ |v'_2|$$

- Sulzmann: ... Let's assume  $v$  is not a *POSIX* value, then there must be another one ... contradiction.
- Exists ?

$$L(r) \neq \emptyset \Rightarrow \exists v. \text{POSIX}(v, r)$$

- In the sequence case  $\text{Seq}(v_1, v_2) \succ_{r_1 \cdot r_2} \text{Seq}(v'_1, v'_2)$ , the induction hypotheses require  $|v_1| = |v'_1|$  and  $|v_2| = |v'_2|$ , but you only know

$$|v_1| @ |v_2| = |v'_1| @ |v'_2|$$

- Although one begins with the assumption that the two values have the same flattening, this cannot be maintained as one descends into the induction (alternative, sequence)

# Our Solution

- A direct definition of what a POSIX value is, using the relation  $s \in r \rightarrow v$  (our specification)

$$\overline{[] \in \epsilon \rightarrow \text{Empty}}$$

$$\overline{[c] \in c \rightarrow \text{Char}(c)}$$

$$\frac{s \in r_1 \rightarrow v}{s \in r_1 + r_2 \rightarrow \text{Left}(v)}$$

$$\frac{s \in r_2 \rightarrow v \quad s \notin L(r_1)}{s \in r_1 + r_2 \rightarrow \text{Right}(v)}$$

$$s_1 \in r_1 \rightarrow v_1$$

$$s_2 \in r_2 \rightarrow v_2$$

$$\neg(\exists s_3 s_4. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r_1) \wedge s_4 \in L(r_2))$$

$$\frac{}{s_1 @ s_2 \in r_1 \cdot r_2 \rightarrow \text{Seq}(v_1, v_2)}$$

...

# Properties

It is almost trivial to prove:

- Uniqueness

If  $s \in r \rightarrow v_1$  and  $s \in r \rightarrow v_2$  then  $v_1 = v_2$

- Correctness

$lexer(r, s) = v$  if and only if  $s \in r \rightarrow v$

# Properties

It is almost trivial to prove:

- Uniqueness

If  $s \in r \rightarrow v_1$  and  $s \in r \rightarrow v_2$  then  $v_1 = v_2$

- Correctness

$lexer(r, s) = v$  if and only if  $s \in r \rightarrow v$

You can now start to implement optimisations and derive correctness proofs for them. But we still do not know whether

$$s \in r \rightarrow v$$

is a POSIX value according to Sulzmann & Lu's definition (biggest value for  $s$  and  $r$ )