

POSIX Lexing with Derivatives of Regular Expressions (Proof Pearl)

Fahad Ausaf, Roy Dyckhoff and Christian Urban

King's College London, University of St Andrews

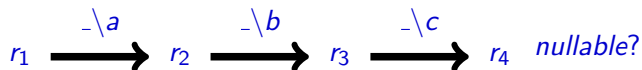
Brzowski's Derivatives of Regular Expressions

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:	$r \setminus []$	$\stackrel{\text{def}}{=} r$
	$r \setminus c::s$	$\stackrel{\text{def}}{=} (r \setminus c) \setminus s$

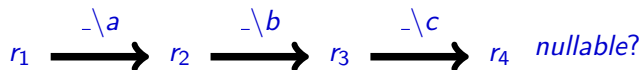
Brzozowski's Matcher

Does r_1 match string abc ?



Brzozowski's Matcher

Does r_1 match string abc ?

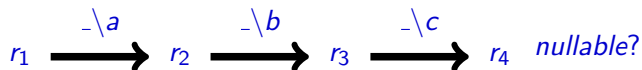


It leads to an elegant functional program:

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

Brzozowski's Matcher

Does r_1 match string abc ?



It leads to an elegant functional program:

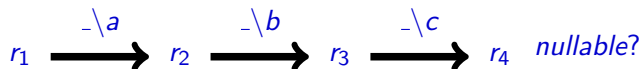
$$\text{matches}(r, s) \stackrel{\text{def}}{=} \text{nullable}(r \backslash s)$$

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

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

Brzowski's Matcher

Does r_1 match string abc ?



It leads to an elegant functional program:

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

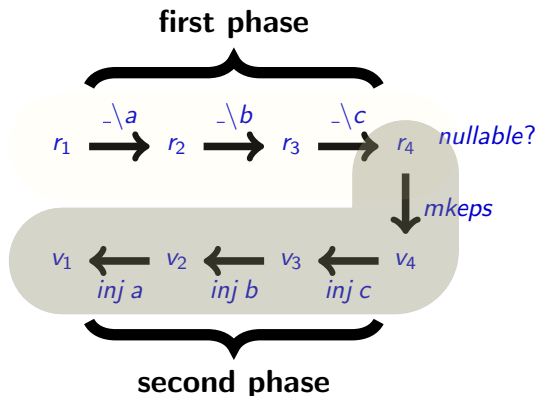
It is an easy exercise to formally prove (e.g. Coq, HOL, Isabelle):

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

But Brzowski'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, ...

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: $r_{keywords} + r_{identifiers}$ (fix graphics below)

```
iffoo_bla
```

```
if_bla
```


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 mathcers are buggy.

http://www.haskell.org/haskellwiki/Regex_Posix

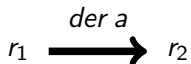
Regular Expressions and Values

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

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

There is also a notion of a string behind a value $|v|$

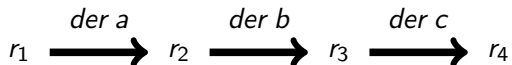
We want to match the string *abc* using r_1



We want to match the string *abc* using r_1



We want to match the string *abc* using r_1

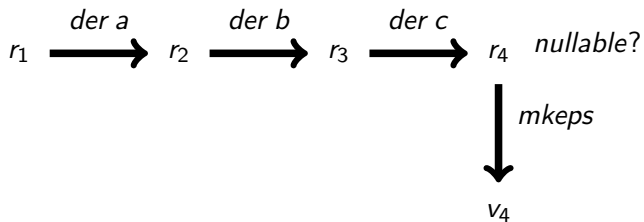


We want to match the string *abc* using r_1

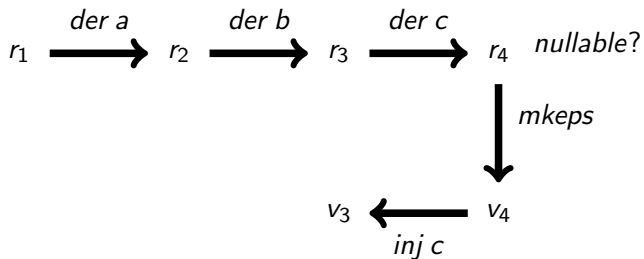


Sulzmann and Lu Matcher

We want to match the string *abc* using r_1

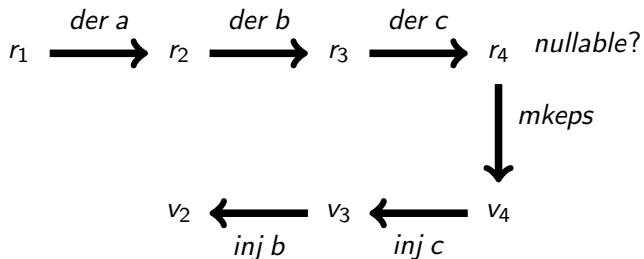


We want to match the string *abc* using r_1

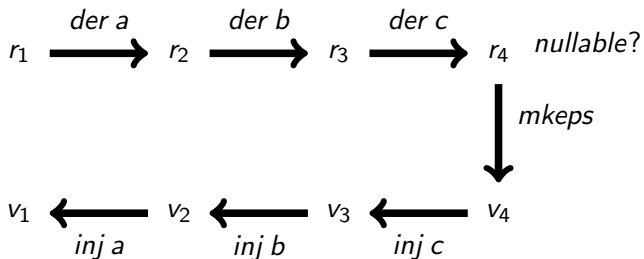


Sulzmann and Lu Matcher

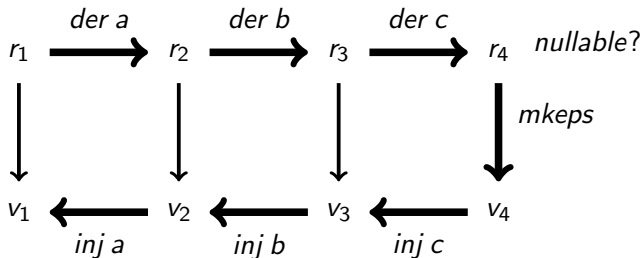
We want to match the string *abc* using r_1



We want to match the string *abc* using r_1



We want to match the string *abc* using r_1



Nullable, Mkeps and Injection Functions

Nullable Function

$nullable(\mathbf{0}) \stackrel{\text{def}}{=} False$

$nullable(\mathbf{1}) \stackrel{\text{def}}{=} True$

$nullable(c) \stackrel{\text{def}}{=} False$

$nullable(r_1 + r_2) \stackrel{\text{def}}{=} nullable\ r_1 \vee nullable\ r_2$

$nullable(r_1 \cdot r_2) \stackrel{\text{def}}{=} nullable\ r_1 \wedge nullable\ r_2$

$nullable(r^*) \stackrel{\text{def}}{=} True$

Mkeps Function

$mkeps(\mathbf{1}) \stackrel{\text{def}}{=} ()$

$mkeps(r_1 \cdot r_2) \stackrel{\text{def}}{=} Seq\ (mkeps\ r_1)\ (mkeps\ r_2)$

$mkeps(r_1 + r_2) \stackrel{\text{def}}{=} \text{if } nullable\ r_1 \text{ then Left } (mkeps\ r_1) \\ \text{else Right } (mkeps\ r_2)$

$mkeps(r^*) \stackrel{\text{def}}{=} Stars\ []$

Injection Function

$inj\ d\ c\ ()$

$inj\ (r_1 + r_2)\ c\ (Left\ v_1)$

$inj\ (r_1 + r_2)\ c\ (Right\ v_2)$

$inj\ (r_1 \cdot r_2)\ c\ (Seq\ v_1\ v_2)$

$inj\ (r_1 \cdot r_2)\ c\ (Left\ (Seq\ v_1\ v_2))$

$inj\ (r_1 \cdot r_2)\ c\ (Right\ v_2)$

$inj\ (r^*)\ c\ (Seq\ v\ (Stars\ vs))$

$\stackrel{\text{def}}{=} Char\ d$

$\stackrel{\text{def}}{=} Left\ (inj\ r_1\ c\ v_1)$

$\stackrel{\text{def}}{=} Right\ (inj\ r_2\ c\ v_2)$

$\stackrel{\text{def}}{=} Seq\ (inj\ r_1\ c\ v_1)\ v_2$

$\stackrel{\text{def}}{=} Seq\ (inj\ r_1\ c\ v_1)\ v_2$

$\stackrel{\text{def}}{=} Seq\ (mkeps\ r_1)\ (inj\ r_2\ c\ v_2)$

$\stackrel{\text{def}}{=} Stars\ (inj\ r\ c\ v::vs)$

- Introduce an inductive defined ordering relation $v \succ_r v'$ which captures the idea of POSIX matching.
- The algorithm returns the maximum of all possible values that are possible for a regular expression.
- The idea is from a paper by Frisch & Cardelli about GREEDY matching (GREEDY = preferring instant gratification to delayed repletion):
- e.g. given $(a + (b + ab))^*$ and string ab

GREEDY: $[Left(a), Right(Left(b))]$

POSIX: $[Right(Right(Seq(a, b)))]$

POSIX Ordering Relation by Sulzmann & Lu

$$\overline{\vdash \text{Empty} : \epsilon}$$

$$\overline{\vdash \text{Char}(c) : c}$$

$$\frac{\vdash v_1 : r_1 \quad \vdash v_2 : r_2}{\vdash \text{Seq}(v_1, v_2) : r_1 \cdot r_2}$$

$$\frac{\vdash v : r_1}{\vdash \text{Left}(v) : r_1 + r_2}$$

$$\frac{\vdash v : r_2}{\vdash \text{Right}(v) : r_1 + r_2}$$

$$\overline{\vdash [] : r^*}$$

$$\frac{\vdash v_1 : r \quad \dots \quad \vdash v_n : r}{\vdash [v_1, \dots, v_n] : r^*}$$

- 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

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

- 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)

- I have no doubt the algorithm is correct — the problem is I do not believe their proof.

“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.”

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

$$\overline{\square} \in \epsilon \rightarrow \textit{Empty}$$

$$\overline{c} \in c \rightarrow \textit{Char}(c)$$

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

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

$$s_1 \in r_1 \rightarrow v_1$$

$$s_2 \in r_2 \rightarrow v_2$$

$$\overline{\neg(\exists s_3 s_4. s_3 \neq \square \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r_1) \wedge s_4 \in L(r_2))}$$

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

...

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)

Conclusions

- We replaced the POSIX definition of Sulzmann & Lu by a new definition (ours is inspired by work of Vansummeren, 2006)
- Their proof contained small gaps (acknowledged) but had also fundamental flaws
- Now, its a nice exercise for theorem proving
- Some optimisations need to be applied to the algorithm in order to become fast enough
- Can be used for lexing, is a small beautiful functional program