

POSIX Lexing with Derivatives of Regular Expressions (Proof Pearl)

Fahad Ausaf
King's College London

joint work with Roy Dyckhoff and Christian Urban

Regular Expressions

$r ::= \emptyset$	null
ϵ	empty string
c	character
$r_1 \cdot r_2$	sequence
$r_1 + r_2$	alternative / choice
r^*	star (zero or more)

The Derivative of a Rexp

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

$der\ cr$ gives the answer, Brzozowski (1964), Owens (2005)
“...have been lost in the sands of time...”

Correctness

It is a relative easy exercise in a theorem prover:

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

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

POSIX Regex Matching

Two rules:

- Longest match rule (“maximal munch rule”): The longest initial substring matched by any regular expression is taken as the next token.

i f f o o _ b l a

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

i f _ b l a

POSIX Regex Matching

Two rules:

- Longest match rule (“maximal munch rule”): The longest initial substring matched by any regular expression is taken as the next token.

i f f o o _ b l a

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

i f _ b l a

Kuklewicz: most POSIX matchers are buggy
http://www.haskell.org/haskellwiki/Regex_Posix

POSIX Regex Matching

- Sulzmann & Lu came up with a beautiful idea for how to extend the simple regular expression matcher to POSIX matching/lexing (FLOPS 2014)



Martin Sulzmann

- the idea: define an inverse operation to the derivatives

Regexes and Values

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

$r ::= \emptyset$	$v ::=$
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)
	$[]$
	$[v_1, \dots, v_n]$

Regexes and Values

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

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

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

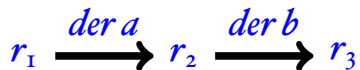
Sulzmann & Lu Matcher

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

$$r_1 \xrightarrow{\text{der } a} r_2$$

Sulzmann & Lu Matcher

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



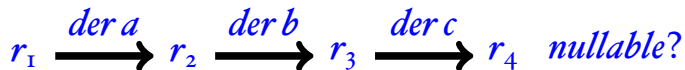
Sulzmann & Lu Matcher

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



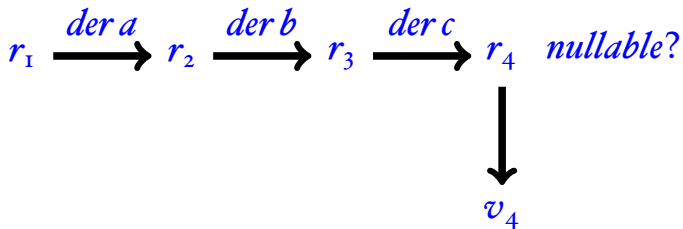
Sulzmann & Lu Matcher

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



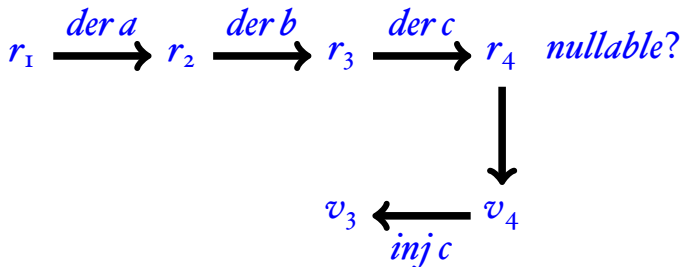
Sulzmann & Lu Matcher

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



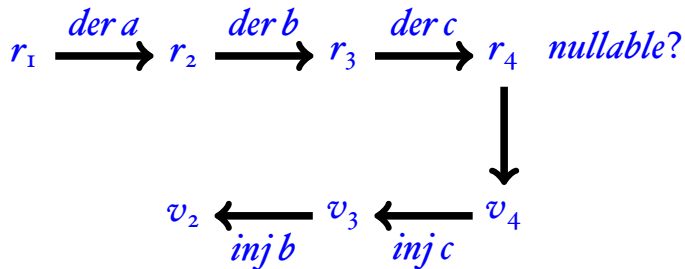
Sulzmann & Lu Matcher

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



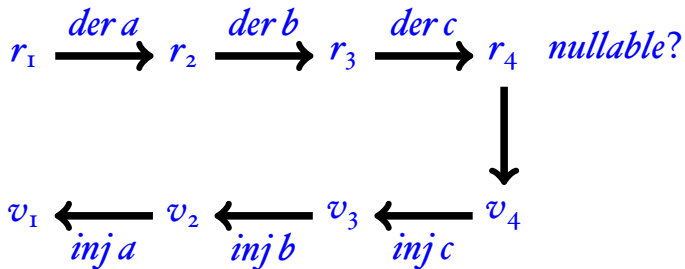
Sulzmann & Lu Matcher

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



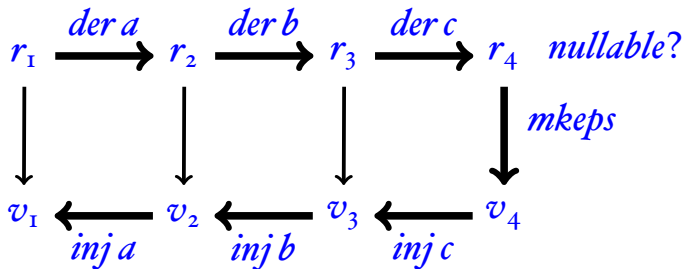
Sulzmann & Lu Matcher

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



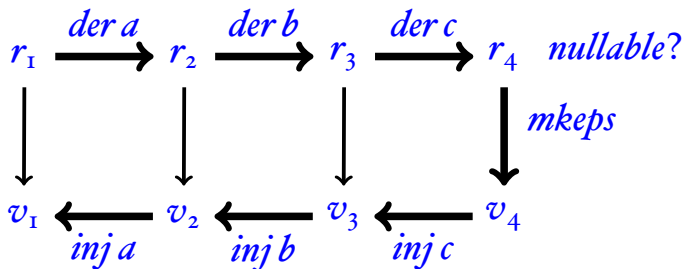
Sulzmann & Lu Matcher

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



Sulzmann & Lu Matcher

We want to match the string abc using r_I :



The original ideas of Sulzmann and Lu are the $mkeps$ and inj functions (omitted here).

Sulzmann & Lu Paper

- 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.”

Sulzmann & Lu Paper

- 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.”

Sulzmann & Lu Paper

- 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

Lemma 3 (Projection and Injection). *Let r be a regular expression, l a letter and v a parse tree.*

- If $\vdash v : r$ and $|v| = lw$ for some word w , then $\vdash \text{proj}_{(r,l)} v : r \setminus l$.*
- If $\vdash v : r \setminus l$ then $(\text{proj}_{(r,l)} \circ \text{inj}_{r \setminus l}) v = v$.*
- If $\vdash v : r$ and $|v| = lw$ for some word w , then $(\text{inj}_{r \setminus l} \circ \text{proj}_{(r,l)}) v = v$.*

MS:BUG[Come across this issue when going back to our constructive reg-ex work] Consider $\vdash [\text{Right } (), \text{Left } a] : (a + \epsilon)^*$. However, $\text{proj}_{((a+\epsilon)^*, a)} [\text{Right } (), \text{Left } a]$ fails! The point is that proj only works correctly if applied on POSIX parse trees.

MS:Possible fixes We only ever apply proj on Posix parse trees.

For convenience, we write “ $\vdash v : r$ is POSIX” where we mean that $\vdash v : r$ holds and v is the POSIX parse tree of r for word $|v|$.

Lemma 2 follows from the following statement.

necessarily hold.

The Proof Idea

by Sulzmann & Lu

- introduce an inductively 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 Proof Idea by Sulzmann & Lu

- introduce an inductively 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 Cardelli & Frisch 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: $[Right(Right(Seq(a, b)))]$

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

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

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

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

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

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

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

Problems

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

Problems

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

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

Problems

- 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

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

Problems

- 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 $s \in r \rightarrow v$ (specification):

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

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

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

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

$$s_1 \in r_1 \rightarrow v_1$$

$$s_2 \in r_2 \rightarrow v_2$$

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

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

...

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)

Conclusion

- 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

Questions?