
Matching regular expressions, revisited, (revisited)*

James McKinna, University of St. Andrews

`james.mckinna@st-andrews.ac.uk`

FUN in Nottingham, February 23, 2007

Last minute GADT hacking: Edwin Brady

Thanks to: Conor, Graham, and Jeremy

The *ur*-problem

Prove that the language recognition problem for regular grammars is

decidable

This should work...

You want decidable? Write a boolean valued function!

```
/* accept : RegExp -> String -> bool */

fun accept r s = acccps r (String.explode s) (fn Nil => true | _ => false)

/* acccps : RegExp -> char list -> (char list -> bool) -> bool */

fun acccps (RE0) cs k = false
  | acccps (RE1) cs k = k cs
  | acccps (REchar a) Nil k = false
  | acccps (REchar a) (c :: cs) k = if (a = c) then k cs else false
  | acccps (REplus r1 r2) cs k = acccps r1 cs k orelse (acccps r2 k cs)
  | acccps (REtimes r1 r2) cs k = acccps r1 cs (fn cs' => acccps r2 cs' k)
  | acccps (r as (REstar r1)) cs k = k cs orelse acccps r1 cs (fn cs' => acccps r cs' k)
```

Problem This function is not *total*. Why?

The *ur*-specification

Imagine a programming language whose type system is expressive enough to allow one to write down, as a *datatype*, the type of:

“successful partial matches of a string *str* against a regexp *r*, exhibited as a parse tree $p : L(r)$, plus a suffix string *s*, such that $str = \text{render } p \ s$ ”

where $\text{render } p \ s$ is the obvious flattening observation on the ADT of parse trees *p*, displaying them as strings.

String recognition with regexps (sketch): typing the recogniser

Consider the inductive family $\text{Recog } r \ s$ of *partially successful matches* with constructor

$$\frac{p : L(r) \quad s : \text{String } A}{\text{consume } p \ s : \text{Recog } r \ (\text{render } p \ s)}$$

which captures successful parsing of *an initial segment* of the input string.

Then a plausible typing for a recogniser is

$$\frac{p : L(r) \quad s : \text{String } A}{\text{recog } p \ s : \text{Recog } r \ s}$$

I: Introduction

Harper (JFP, 1999): “proof-directed debugging”

A Lakatos-style “failed proof” analysis, with repair, of:

- a continuation-passing style regexp matcher in SML
- partial correctness: tricky informal inductive argument, *assuming* inductively that continuations involved terminate: “*validity*”
- termination: *fails* for matching against 0^*
- repair:
 - eliminate null transitions: compute $\delta(r) = 0, 1$
 - put r in *standard form*: $r = \delta(r) + r^-$
- termination secured for $0, 1$ and regexps of the form r^-

Harper's broken matcher

```
/* accept : RegExp -> String -> bool */

fun accept r s = acccps r (String.explode s) (fn Nil => true | _ => false)

/* acccps : RegExp -> char list -> (char list -> bool) -> bool */

fun acccps (RE0) cs k = false
| acccps (RE1) cs k = k cs
| acccps (REchar a) Nil k = false
| acccps (REchar a) (c :: cs) k = if (a = c) then k cs else false
| acccps (REplus r1 r2) cs k = acccps r1 cs k orelse (acccps r2 k cs)
| acccps (REtimes r1 r2) cs k = acccps r1 cs (fn cs' => acccps r2 cs' k)
| acccps (r as (REstar r1)) cs k = k cs orelse acccps r1 cs (fn cs' => acccps r cs' k)
```

This function is not structural recursive because of the highlighted recursive call on `acccps`, and there is no guarantee that any of the string will have been consumed at that point.

A natural candidate for formalisation?

- Tricky inductive correctness argument
- Need to analyse possible non-termination of SML programs
- Re-frame the problem!
 - Work in a theory of *terminating* functions (EPIGRAM)
 - Proof by induction, definition by structural recursion
 - Distinguish *types* of regexps, resp. those in standard form; a *view* (Wadler (1987); McBride/McKinna (2004))
 - Matching can itself also be expressed as a *view* (of Strings)
 - Correctness evident by *type* (Curry-Howard as usual)
 - Lose CPS; recover direct-style matcher (continuations in the tail)

end of part I

II: Regular expressions and regular languages

Why dependent types matter

Our treatment makes use of dependent types to represent, in a uniform framework,

- the language sets $L(r)$ for a given regexp r ;
- the function rendering each $p : L(r)$ as a string;
- derivation trees $i : r \leq s$ axiomatising the inclusion $L(r) \subseteq L(s)$;
- the mapping $\llbracket i \rrbracket : L(r) \rightarrow L(s)$, together with $\langle\langle i \rangle\rangle$ (which witnesses this inclusion), the proof that $\llbracket i \rrbracket$ preserves renderings;

in other words, the syntax, semantics *and* proof theory of regular languages.

The type of regular expressions

RE A ,

- defined over an alphabet A ,
- is definable as usual, with
 - 0 ,
 - 1 ,
 - characters a ,
 - alternation $r + s$,
 - composition $r . s$,
 - and repetition (Kleene star) r^* ; iteration, $r^+ \equiv r . r^*$

The inductive family of language sets

We do *not* define $L(r)$ directly as sets of strings, but as abstract syntax trees for successful parses of such strings:

data $\frac{r : \text{RE } A}{L(r) : \star}$ where \dots

where $\frac{}{\epsilon : L(\mathbf{1})}$ $\frac{a : A}{[a] : L(a)}$ $\frac{p : L(r) ; q : L(s)}{p \bullet q : L(r \cdot s)}$

$\frac{p : L(r)}{l(p) : L(r + s)}$ $\frac{q : L(s)}{r(q) : L(r + s)}$

$\frac{}{\epsilon : L(r^*)}$ $\frac{p : L(r) ; q : L(r^*)}{p \star q : L(r^*)}$

NB. $L(\mathbf{0})$ is an *empty* type: no constructor is declared for it.

Alternatively: a recursive definition

$$\underline{\text{let}} \quad \frac{r : \text{RE } A}{L(r) : \star}$$

$$L(\mathbf{0}) \quad \Rightarrow \quad \mathbf{0}$$

$$L(\mathbf{1}) \quad \Rightarrow \quad \mathbf{1}$$

$$L(a) \quad \Rightarrow \quad \{a\}$$

$$L(r \cdot s) \quad \Rightarrow \quad L(r) \times L(s)$$

$$L(r + s) \quad \Rightarrow \quad L(r) \oplus L(s)$$

$$L(r^*) \quad \Rightarrow \quad \text{List } (L(r))$$

Rendering parse trees as strings: obvious fusion

$$\underline{\text{let}} \quad \frac{p : L(r) ; str : \text{String } A}{\text{render } p \text{ str} : \text{String } A}$$

$$\text{render } \epsilon \text{ str} \quad \Rightarrow \quad str$$

$$\text{render } [a] \text{ str} \quad \Rightarrow \quad \text{prefix } a \text{ str}$$

$$\text{render } (p \bullet q) \text{ str} \Rightarrow \text{render } p (\text{render } q \text{ str})$$

$$\text{render } (l(p)) \text{ str} \Rightarrow \text{render } p \text{ str}$$

$$\text{render } (r(q)) \text{ str} \Rightarrow \text{render } q \text{ str}$$

$$\text{render } \varepsilon \text{ str} \quad \Rightarrow \quad str$$

$$\text{render } (p \star q) \text{ str} \Rightarrow \text{render } p (\text{render } q \text{ str})$$

An obvious 'fusion' lemma

exercise: prove the following

$$\text{render } p \text{ } str = \text{append} (\text{render } p \text{ []}) str$$

with [] the null string, and string concatenation **append** as usual.

Standard form regular expressions

The crux of the termination argument:

- identify the sublanguage of expressions $s : \text{SRE } A$
- such that their parse trees $p : \mathbf{L}(s)$ render as *non-empty* strings.
- (matching against such an s is guaranteed to terminate, by well-founded induction on string suffixes, using the fusion lemma)

$$\text{data } \frac{A : \star}{\text{SRE } A : \star} \quad \text{where } \frac{}{0 : \text{SRE } A} \quad \frac{a : A}{a : \text{SRE } A}$$

$$\dots \quad \frac{s, s' : \text{SRE } A}{s + s' : \text{SRE } A} \quad \frac{s, s' : \text{SRE } A}{s . s' : \text{SRE } A} \quad \frac{s : \text{SRE } A}{s^+ : \text{SRE } A}$$

Standardising regular expressions: the theorem

- An obvious erasure from **SRE** A to **RE** A , $|s|$
- re-presentation of Harper's analysis by proving the following normal form theorem for regexps: every $r : \text{RE } A$ is equivalent to either
 - $|s|$ for $s : \text{SRE } A$, or
 - $1 + |s|$ for $s : \text{SRE } A$
- 'equivalent' here means: generates the same strings via **render**

Standardising regular expressions: the function

- the two cases of the theorem correspond to two constructors of a type family, $\text{Std } r$, (the cases where $\delta(r) = 0$, $\delta(r) = 1$ in Harper)

- data $\frac{r : \text{RE } A}{\text{Std } r : \star}$ where

$$\frac{s : \text{SRE } A \quad i : |s| \equiv r}{\text{Std}_0 s : \text{Std } |s|} \quad \frac{s : \text{SRE } A \quad i : 1 + |s| \equiv r}{\text{Std}_1 s i : \text{Std } r}$$

- the theorem is witnessed by a function of type

$$\frac{r : \text{RE } A}{\text{std } r : \text{Std } r}$$

- incl. all the machinery behind $r \equiv \delta(r) + r^-$

Axiomatising the equational theory of regular expressions

- $r \equiv s$ is (a fragment of) the familiar equational theory of Kleene algebra (semantically, “ $\mathbf{L}(r) \equiv \mathbf{L}(s)$ ”)
- technically easier to mix the equational and inequational theory, capturing “ $r \leq r'$ ” as well: just another inductive family...
- crucial soundness lemma: if $r \leq r'$, then $\mathbf{L}(r) \subseteq \mathbf{L}(r')$
- this lemma becomes a *function*, factorised as
 - a mapping $\llbracket i \rrbracket : \mathbf{L}(r) \rightarrow \mathbf{L}(r')$; defined by structural induction over $i : r \leq r'$, respectively $i : r \equiv r'$;
 - a proof $\langle\langle i \rangle\rangle$ that $\mathbf{render} (\llbracket i \rrbracket p) \text{ str} = \mathbf{render} p \text{ str}$;

end of part II

III: Writing the matcher

Specifying the matcher

- The problem of (partial) matching a given string str against a regexp r is then to construct
 - a parse tree $p : L(r)$ and
 - a suffix sfx
 - such that $str = \text{render } p \text{ } sfx$,that is, to invert the rendering function.
- The classical recognition problem $str? \in L(r)$ then reduces to the problem of testing whether the suffix $sfx = []$

Specifying the family

$$\text{data} \quad \frac{r : \text{RE } A ; \text{str} : \text{String } A}{\text{Recog } r \text{ str} : \star}$$

$$\text{where} \quad \frac{p : \text{L}(r) ; \text{sfx} : \text{String } A}{\text{consume } p \text{ sfx} : \text{Recog } r (\text{render } p \text{ sfx})}$$

$$\frac{e : \text{Err } r}{\text{err } e : \text{Recog } r \bar{e}}$$

- Writing a recogniser then amounts to writing a function **recog** of type

$$\forall r : \text{RE } A. \forall s : \text{String } A. \text{Recog } r s$$

- Non-dependent elimination (the technique of **views**) over the family **Recog** $r s$ exposes s so as to invert **render**

CPS revisited: why partial matching matters

- when attempting to match str against a composition $r . s$, we can recursively match the suffix sfx against s
- what holds this together is that the corresponding result types match up, because for $p : L(r)$, $q : L(s)$ we have as a *definitional* equality

$$\text{render } (p \bullet q) \text{ sfx} = \text{render } p \text{ (render } q \text{ sfx)}.$$

- similar considerations apply when matching against r^+ etc.
- as in Wand's influential analysis, the suffix strings encode, as usual, the continuation of the computation

The recogniser

- We define the recogniser as a function `recog`, declared with signature

$$\underline{\text{let}} \quad \frac{r : \text{RE } A ; \text{str} : \text{String } A}{\text{recog } r \text{ str} : \text{Recog } r \text{ str}}$$

- `recog r str` is computed as follows:

- standardise r , yielding $s : \text{SRE } A$, a constructor tag indicating whether r recognises ϵ , and a proof $i : s \leq r$;

- match using a specialised recogniser `srecog`, declared with

signature $\underline{\text{let}} \quad \frac{s : \text{SRE } A \quad \text{str} : \text{String } A}{\text{srecog } s \text{ str} : \text{Recog } |s| \text{ str}}$

- use the proof $\langle\langle i \rangle\rangle$ to fix up the types!

The specialised recogniser

Matching a specialised recogniser `srecog`, declared with signature

let $\frac{s : \text{SRE } A \quad str : \text{String } A}{\text{srecog } s \text{ } str : \text{Recog } |s| \text{ } str}$ does the obvious thing:

- fail on `0`
- on `a`: succeed if the head character is matched; fail otherwise
- on `+`: try to match the left, try the right if you fail
- on `.`: try to match an initial segment, continue with the tail
- on `s+`: try to match one copy of `s`... once you fail, return the suffix and figure out if you have succeeded or failed!

Termination relies on the fact that you must consume tokens at each success step which gives rise to a recursive call

end of part III

Conclusions

- relativisation to a given r makes $L(r)$, and **render**, evidently “correct”
- (category-theoretic) treatment of Kleene algebra: $L()$ is a *functor*
- Harper-Sethi/Berry-McNaughton/Yamada normalisation to **SRE** A
- well-founded recursion on suffices secures termination for regexps in standard form
- the **consume** constructor encodes, direct-style, (the continuation on) the suffix of the string, having successfully parsed a prefix of the string
- the matcher is itself, evidently “correct” by virtue of its type

Anti-Conclusions

- the eventual program
 - has a very high ‘deBruijn ratio’ compared to Harper’s original
 - still isn’t entirely finished. . . oops!
 - is harder to understand, . . . or is it?
- OTT might help hide all the $\langle\langle i \rangle\rangle$ equational reasoning in types

Questions?

Dependent families of types [Martin-Löf 1971]

- The key device we exploit to achieve this is the idea of a

dependent family of types $F : T \rightarrow \star$

a function on type $T : \star$ which returns *types* $F t : \star$ given $t : T$.

- allow *arbitrary* T as the domain of variation (not just \star itself)
- then F behaves like a *predicate* on T
- quantification \forall, \exists given by type constructors $\Pi, \Sigma \dots$ so we have typed programs and logic with explicit proofs
- an important class of datatypes arise by considering *inductively-defined* F [Dybjer 1991].

A uniform generalisation of GADTs and related notions

the spectrum of possible instances $T \vec{a}$ occurring in source and target types of term constructors and functions:

Hindley-Milner \vec{a} can be type variables only; *uniform* choice over all constructors of a datatype; function instances $T \vec{\tau}$ similarly uniform

polymorphic recursion non-uniform instances in *source* types for constructors, on a per-constructor basis

GADTs non-uniform instances in source and *target* types: \vec{a} may be *arbitrary* type expressions (necessarily in type-constructor form; no type-level functions)

Ω mega \vec{a} may be *arbitrary* type expressions (*not* necessarily in type-constructor form)

Inductive Families

- Dybjer's families: \vec{a} may be *term* expressions (*not* necessarily types!)
- can consider further stratifications of this idea
 - only consider \vec{a} to be *variables* (Cayenne)
 - ... to be *constructor form patterns*
 - *arbitrary* expressions
- EPIGRAM makes the last, most permissive choice
- NB. these types are **not** ascribed to untyped terms
- they *prescribe*, and/or *describe* very rich properties
- don't understand data or computations *independently* of their types
- give up all partial recursive functions. . .

Phase distinction: the event horizon for type systems

There seems to be a fundamental problem with keeping static and dynamic layers apart:

you cannot say that the thing you construct is
related to the input you started with

For others, there seems to be a fundamental problem with mixing static and dynamic layers:

types might (have to) get passed at run-time