

Structurally Recursive Descent Parsing

Nils Anders Danielsson (Nottingham)
Joint work with Ulf Norell (Chalmers)

Parser combinators

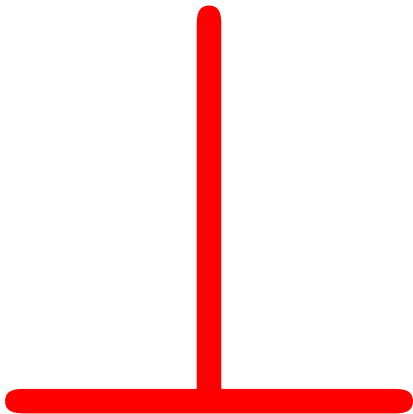
- ▶ Parser combinator libraries are great!
- ▶ Elegant code.
- ▶ Executable grammars.
- ▶ Easy to abstract out recurring patterns.
- ▶ Light-weight.
- ▶ Nowadays often fast enough.

Simple example

$expr = _ + _ \ \$ \ term \ \langle * \ sym \ ' + ' \ \langle * \rangle \ expr$
|
 $term = \dots$
 $term$

Simple example

$expr = _ + _ \quad \$ \quad expr \langle * \quad sym \text{'+' } \langle * \rangle \quad term$
|
 $term = \dots$



Risk of non-termination

- ▶ Combinator parsing is **not guaranteed to terminate**.
- ▶ Most combinator parsers fail for left-recursive grammars.
- ▶ *Executable* grammars?
- ▶ Some errors are not caught at compile-time.

Another problem

- ▶ All interesting grammars are **cyclic**:

$$\begin{array}{l} \text{expr} = \text{--+-} \$ \text{ term} \langle * \text{ sym ' + ' } \langle * \rangle \text{ expr} \\ | \qquad \qquad \qquad \text{term} \end{array}$$

- ▶ Cyclic values are hard to understand and reason about.
- ▶ How do you implement combinator parsing in a language which requires structural recursion?

Our solution

- ▶ Remove cycles by representing grammars as functions from non-terminals to parsers:

Grammar tok nt = nt res \rightarrow *Parser tok nt res*

- ▶ Rule out left recursion by restricting the types.

Examples

Example

- ▶ Non-terminals:

data *NT* : *ParserType* **where**

expr : *NT* _ \mathbb{N}

term : *NT* _ \mathbb{N}

- ▶ Result type: \mathbb{N} .
- ▶ Indices ensuring termination: $_$.
Inferred automatically.

Example

$g : \text{Grammar Char NT}$

$g \text{ expr} = _ + _ \ \$ \ ! \ \text{term} \ \langle * \ \text{sym} \ ' + ' \ \langle * \rangle \ ! \ \text{expr}$
 | !
 $g \text{ term} = \dots$

- ▶ Uses applicative functor interface.
- ▶ Monadic interface also possible.

Abstraction

- ▶ Much of the flavour of ordinary combinator parsers is preserved.
- ▶ Abstraction requires a little work, though.

Abstraction

data *NT* : *ParserType* **where**

lib : *L.Nonterminal NT i r* \rightarrow *NT _ r*

expr : *NT _* \mathbb{N}

term : *NT _* \mathbb{N}

op : *NT _* ($\mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathbb{N}$)

open L.Combinators lib

g : *Grammar Char NT*

g (*lib p*) = *libraryGrammar p*

g expr = *chainl*₁ (! *term*) (! *op*)

g term = *number* | *parenthesised* (! *expr*)

g op = *_+_* <\$ *sym* '+'

| *_--* <\$ *sym* '-'

Abstraction

```
data NT : ParserType where  
  lib   : L.Nonterminal NT i r → NT _ r  
  expr  : NT _ ℕ  
  term  : NT _ ℕ  
  op    : NT _ (ℕ → ℕ → ℕ)  
open L.Combinators lib  
g : Grammar Char NT  
g (lib p) = libraryGrammar p  
g expr   = chainl1 (! term) (! op)  
g term   = number | parenthesised (! expr)  
g op     = _+_ <$ sym '+'  
          | _-_- <$ sym '-'
```

Running a parser

parse : *Parser tok nt i result*
→ *Grammar tok nt*
→ [*tok*] → [*result* × [*tok*]]

How does it work?

Indices

Parsers are indexed on two things:

$$\text{Index} = \text{Empty} \times \text{Corners}$$

Empty Does the parser accept the empty string?

Corners A tree representation of the proper left corners of the parser.

Indices

Empty Does the parser accept the empty string?

Empty = Bool

Corners Represents all positions in the grammar in which the parser must not recurse to itself.

data *Corners* : *Set* **where**

leaf : *Corners*

step : *Corners* \rightarrow *Corners*

node : *Corners* \rightarrow *Corners* \rightarrow *Corners*

Some basic combinators

- <*>- : $\text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow$
 $\text{Parser } (e_1 \wedge e_2, \text{if } e_1 \text{ then node } c_1 \text{ } c_2 \text{ else } c_1)$
- |- : $\text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow$
 $\text{Parser } (e_1 \vee e_2, \text{node } c_1 \text{ } c_2)$
- !_ : $nt (e, c) \rightarrow \text{Parser } (e, \text{step } c)$

Some basic combinators

- $_ \langle * \rangle _$: $\text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow$
 $\text{Parser } (e_1 \wedge e_2, \text{if } e_1 \text{ then node } c_1 \ c_2 \text{ else } c_1)$
- $_ | _$: $\text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow$
 $\text{Parser } (e_1 \vee e_2, \text{node } c_1 \ c_2)$
- $! _$: $nt (e, c) \rightarrow \text{Parser } (e, \text{step } c)$

Some basic combinators

- $_ \langle * \rangle _ : \text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow$
 $\text{Parser } (e_1 \wedge e_2, \text{if } e_1 \text{ then node } c_1 \text{ } c_2 \text{ else } c_1)$
- $_ | _ : \text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow$
 $\text{Parser } (e_1 \vee e_2, \text{node } c_1 \text{ } c_2)$
- $! _ : nt (e, c) \rightarrow \text{Parser } (e, \text{step } c)$

This does not type check:

grammar : nt (e, c) res \rightarrow Parser tok nt (e, c) res
grammar rec = ! rec

Reason: $c \neq \text{step } c$.

Some basic combinators

- $_ \langle * \rangle _ : \text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow \text{Parser } (e_1 \wedge e_2, \text{if } e_1 \text{ then node } c_1 \text{ } c_2 \text{ else } c_1)$
- $_ | _ : \text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow \text{Parser } (e_1 \vee e_2, \text{node } c_1 \text{ } c_2)$
- $! _ : nt (e, c) \rightarrow \text{Parser } (e, \text{step } c)$

This works, though:

grammar : nt (e, c) res → Parser tok nt (e, c) res
*grammar rec = sym c * > ! rec*

Reason: *sym c* must consume a token.

Some basic combinators

- $_ \langle * \rangle _ : \text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow \text{Parser } (e_1 \wedge e_2, \text{if } e_1 \text{ then node } c_1 \text{ } c_2 \text{ else } c_1)$
- $_ | _ : \text{Parser } (e_1, c_1) \rightarrow \text{Parser } (e_2, c_2) \rightarrow \text{Parser } (e_1 \vee e_2, \text{node } c_1 \text{ } c_2)$
- $! _ : nt (e, c) \rightarrow \text{Parser } (e, \text{step } c)$

Indirect left recursion also fails:

grammar : $nt (e, c) \text{ res} \rightarrow \text{Parser } tok \text{ } nt (e, c) \text{ res}$
grammar $rec = ! \text{ other } \langle * \rangle \dots$
grammar $\text{other} = \text{many } p \langle * \rangle ! \text{rec}$

Indices can be useful anyway

- ▶ Parsing zero or more things:

$$\begin{aligned} \textit{many} & : \textit{Parser tok nt} (\textit{false}, c) r \\ & \rightarrow \textit{Parser tok nt} _ \quad [r] \end{aligned}$$

- ▶ Note that the input parser must not accept the empty string.
- ▶ Even if the backend can handle *many empty* it seems reasonable to assume that it is a bug.

Backend

- ▶ Simple backtracking implementation. (So far.)
- ▶ Lexicographic structural recursion over:
 1. An upper bound on the length of the input string.
 2. The *Corners* index.
 3. The structure of the parser.

Expressive power?

- ▶ Can define grammars with an infinite number of non-terminals:

data *NT* : *ParserType* **where**

$a^{1+} _ : \mathbb{N} \rightarrow NT _ Unit$

g : *Grammar Char NT*

$g (a^{1+} \text{ zero}) = sym 'a' * > return unit$

$g (a^{1+} (\text{suc } n)) = sym 'a' * > !(a^{1+} n)$

- ▶ Can use this to define non-context-free languages: $a^n b^n c^n$.

Expressive power?

- ▶ Can define grammars with an infinite number of non-terminals:

data *NT* : *ParserType* **where**

$a^{1+} _ : \mathbb{N} \rightarrow NT _ Unit$

g : *Grammar Char NT*

$g (a^{1+} \text{ zero}) = sym 'a' * > return unit$

$g (a^{1+} (\text{suc } n)) = sym 'a' * > !(a^{1+} n)$

- ▶ **Careful!** Types can become really complicated:

$nt : (n : \mathbb{N}) \rightarrow NT (f n) Unit$

Expressive power?

- ▶ Can define grammars with an infinite number of non-terminals:

data *NT* : *ParserType* **where**

$a^{1+} _ : \mathbb{N} \rightarrow NT _ Unit$

g : *Grammar Char NT*

$g (a^{1+} \text{ zero}) = sym 'a' * > return unit$

$g (a^{1+} (\text{suc } n)) = sym 'a' * > !(a^{1+} n)$

- ▶ The same warning applies when defining libraries.

Conclusions

- ▶ Structurally recursive descent parsing.
- ▶ Termination guaranteed.
- ▶ Errors caught at compile-time.
- ▶ Still feels like combinator parsing.
- ▶ More complicated types,
but the overhead for the user is usually small.

Possible future work

- ▶ More efficient backend.
- ▶ Use backend which can handle left recursion \Rightarrow less complicated types.
 - ▶ But the types can be nice to have anyway.
 - ▶ And who needs left recursion?
chainl is more high-level.

?

Extra slides

Defining a library: Non-terminals

The non-terminals are parameterised on the outer grammar's non-terminals:

```
data NT (nt : ParserType) : ParserType where  
  many : Parser tok nt (false, c) r → NT nt _ [r]  
  many1 : Parser tok nt (false, c) r → NT nt _ [r]
```

Defining a library: Combinators

Combinators parameterised on a *lib* constructor:

module *Combinators* (*lib* : *NT nt i r* → *nt i r*) **where**

^{*} : *Parser tok nt* (false, *c*) *r* → *Parser tok nt* ** [*r*]

p^{*} = ! *lib* (many *p*)

⁺ : *Parser tok nt* (false, *c*) *r* → *Parser tok nt* ** [*r*]

p⁺ = ! *lib* (many₁ *p*)

library : *NT nt i r* → *Parser tok nt i r*

library (many *p*) = *return* [] | *p*⁺

library (many₁ *p*) = *_*::*_* \$ *p* <*> *p*^{*}

Defining a library: Combinators

Wrappers (to ease use of the library):

module *Combinators* (*lib* : *NT nt i r* \rightarrow *nt i r*) **where**

$_{}^*$: *Parser tok nt* (*false, c*) *r* \rightarrow *Parser tok nt* $_{} [r]$

p^* = ! *lib* (*many p*)

$_{}^+$: *Parser tok nt* (*false, c*) *r* \rightarrow *Parser tok nt* $_{} [r]$

p^+ = ! *lib* (*many₁ p*)

library : *NT nt i r* \rightarrow *Parser tok nt i r*

library (*many p*) = *return []* | p^+

library (*many₁ p*) = $_{}::_{} \$ p <*> p^*$

Defining a library: Combinators

Grammar (as before):

module *Combinators* (*lib* : *NT nt i r* \rightarrow *nt i r*) **where**

$_*$: *Parser tok nt* (*false, c*) *r* \rightarrow *Parser tok nt* $_$ [*r*]

p^* = ! *lib* (*many* *p*)

$_+$: *Parser tok nt* (*false, c*) *r* \rightarrow *Parser tok nt* $_$ [*r*]

p^+ = ! *lib* (*many*₁ *p*)

library : *NT nt i r* \rightarrow *Parser tok nt i r*

library (*many* *p*) = *return* [] | p^+

library (*many*₁ *p*) = $_::_ \$ p <*> p^*$