

Structurally Recursive Descent Parsing

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

Cambridge, 2009-03-20

Parser combinator example

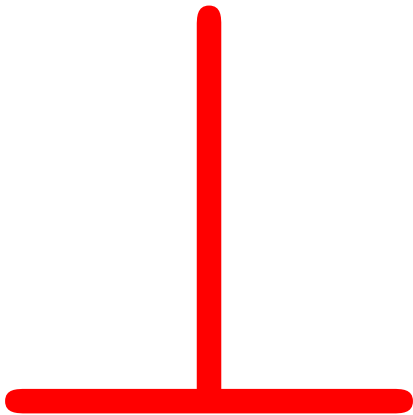
$$\begin{array}{l} \text{expr} = \text{term} \cdot \text{tok "+"} \cdot \text{expr} \\ \quad \quad | \text{term} \\ \text{term} = \dots \end{array}$$

Parser combinator propaganda

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

Parser combinator example

$$\begin{array}{l} \text{expr} = \text{expr} \cdot \text{tok "+"} \cdot \text{term} \\ \quad \quad | \text{ term} \\ \text{term} = \dots \end{array}$$



Risk of non-termination

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

Goal

Parser combinator library implemented in total language

- ▶ For simplicity: recursive descent.

Language: Agda

- ▶ Functional.
- ▶ Dependent types.
- ▶ Total.
- ▶ Support for mixed induction/coinduction and mutual structural recursion/guarded corecursion.

Recognisers

Termination

- ▶ Recursive descent \Rightarrow left recursion has to be avoided.
- ▶ Grammars often cyclic \Rightarrow corecursion.
- ▶ Recognition returns inductive value (*Bool*) \Rightarrow structural recursion.

Interface

\emptyset : G

ε : G

tok : $Token \rightarrow G$

$_|_$: $G \rightarrow G \rightarrow G$

$_._$: $G \rightarrow G \rightarrow G$

Corecursion will be used \Rightarrow
some arguments have to be coinductive.

Choice

Hard to decide infinite choice:

$$p = p \mid p'$$

$$p = p' \mid p$$

The arguments of $_|_$ should be inductive.

Sequencing

Also problematic (due to top-down, left-to-right implementation):

$$p = p \cdot p'$$

OK if p' does not accept the empty string:

$$p = p' \cdot p$$

Grammars

Index **true** iff empty string accepted:

data $G : Bool \rightarrow Set$ **where**

\emptyset : G **false**

ε : G **true**

tok : $Token \rightarrow G$ **false**

$_|_$: $\forall \{e_1 e_2\} \rightarrow G e_1 \rightarrow G e_2 \rightarrow G (e_1 \vee e_2)$

$_?._$: $\forall \{e_2\} \rightarrow G$ **true** $\rightarrow G e_2 \rightarrow G e_2$

$_!._$: $\forall \{e_2\} \rightarrow G$ **false** $\rightarrow \infty (G e_2) \rightarrow G$ **false**

Coinductive types

- ▶ ∞ marks coinductive arguments.
- ▶ Can be seen as a suspension.
- ▶ Delay and force:

$$\begin{aligned} \# &: \forall \{A\} \rightarrow A \rightarrow \infty A \\ \flat &: \forall \{A\} \rightarrow \infty A \rightarrow A \end{aligned}$$

- ▶ $T = F (\infty T) T \approx T = \nu C. \mu l. F C l.$

Coinductive types

$G : Bool \rightarrow Set$

$G e \cong \nu C. \mu l.$

$e \equiv false$

+ $e \equiv true$

+ $Token \times e \equiv false$

+ $\Sigma e_1 e_2. l e_1 \times l e_2 \times e \equiv e_1 \vee e_2$

+ $\Sigma e_2. l true \times l e_2 \times e \equiv e_2$

+ $\Sigma e_2. l false \times C e_2 \times e \equiv false$

Semantics

Inductive:

data $_ \in _$: $\forall \{e\} \rightarrow List\ Token \rightarrow G\ e \rightarrow Set$ **where**

ε : $[\] \in \varepsilon$

tok : $[t] \in \mathbf{tok}\ t$

$|^l$: $s \in p_1 \rightarrow s \in p_1 \mid p_2$

$|^r$: $s \in p_2 \rightarrow s \in p_1 \mid p_2$

$_?._$: $s_1 \in p_1 \rightarrow s_2 \in p_2 \rightarrow s_1 \uparrow\uparrow s_2 \in p_1\ ?.\ p_2$

$_!_$: $s_1 \in p_1 \rightarrow s_2 \in p_2 \rightarrow s_1 \uparrow\uparrow s_2 \in p_1\ !.\ p_2$

For $p : G\ e$:

$[\] \in p$ iff $e \equiv \mathbf{true}$.

Example

Kleene star:

mutual

$_ \star : G \text{ false} \rightarrow G \text{ true}$

$p \star = \varepsilon \mid p \mid p \star$

$_ + : G \text{ false} \rightarrow G \text{ false}$

$p + = p ! \cdot \# (p \star)$

The argument must not accept the empty string;
 $\varepsilon \star$ is not very useful.

Implementation

Implementation

Bounded vectors:

data $BVec$ ($A : Set$) : $\mathbb{N} \rightarrow Set$ **where**

$[]$: $\forall \{n\} \rightarrow BVec A n$

$x :: xs$: $\forall \{n\} \rightarrow A \rightarrow BVec A n \rightarrow BVec A (suc n)$

\uparrow : $\forall \{A n\} \rightarrow BVec A n \rightarrow BVec A (suc n)$

$\uparrow [] = []$

$\uparrow (x :: xs) = x :: \uparrow xs$

Implementation

$parse : \forall \{e\ n\} \rightarrow G\ e \rightarrow BVec\ Token\ n \rightarrow$
 $List\ (BVec\ Token\ (if\ e\ then\ n\ else\ pred\ n))$

$parse\ \emptyset\ s = fail$

$parse\ \epsilon\ s = return\ s$

$parse\ (tok\ t)\ [] = fail$

$parse\ (tok\ t)\ (t' :: s) =$
 $if\ t == t'\ then\ return\ s\ else\ fail$

Implementation

$$\begin{aligned} \text{parse} &: \forall \{e\ n\} \rightarrow G\ e \rightarrow BVec\ Token\ n \rightarrow \\ &\quad List\ (BVec\ Token\ (if\ e\ then\ n\ else\ pred\ n)) \\ \text{parse}\ (p_1\ |_{e_1}\ p_2)\ s &= (\uparrow_1\ e_1\ \langle \$ \rangle\ \text{parse}\ p_1\ s) \text{ ++ } \\ &\quad (\uparrow_2\ e_1\ \langle \$ \rangle\ \text{parse}\ p_2\ s) \\ \text{parse}\ (p_1\ ?\cdot\ p_2)\ s &= \text{parse}\ p_2 \preccurlyeq \text{parse}\ p_1\ s \\ \text{parse}\ (p_1\ !\cdot\ p_2)\ [] &= \text{fail} \\ \text{parse}\ (p_1\ !\cdot_{e_2}\ p_2)\ (t\ ::\ s) &= \\ &\quad \uparrow_3\ e_2\ \langle \$ \rangle\ (\text{parse}\ (b\ p_2) \preccurlyeq \text{parse}\ p_1\ (t\ ::\ s)) \end{aligned}$$

Implementation

$parse : \forall \{e\ n\} \rightarrow G\ e \rightarrow BVec\ Token\ n \rightarrow$
 $List\ (BVec\ Token\ (if\ e\ then\ n\ else\ pred\ n))$

$parse\ (p_1\ |_{e_1}\ p_2)\ s = (\uparrow_1\ e_1\ \langle \$ \rangle\ parse\ p_1\ s) \ ++$
 $(\uparrow_2\ e_1\ \langle \$ \rangle\ parse\ p_2\ s)$

$parse\ (p_1\ ?\cdot\ p_2)\ s = parse\ p_2 \ \lll\ parse\ p_1\ s$

$parse\ (p_1\ !\cdot\ p_2)\ [] = fail$

$parse\ (p_1\ !\cdot_{e_2}\ p_2)\ (t :: s) =$
 $\uparrow_3\ e_2\ \langle \$ \rangle\ (parse\ p_2) \ \lll\ parse\ p_1\ (t :: s)$

Implementation

$_ \in? _ : \forall \{e\} \rightarrow List\ Token \rightarrow G\ e \rightarrow Bool$
 $s \in? p = any\ null\ (parse\ p\ (fromList\ s))$

Parsers

Actual parsers

data $P : Bool \rightarrow Set \rightarrow Set$ **where**

- return** : $\forall \{R\} \rightarrow R \rightarrow P \text{ true } R$
- fail** : $\forall \{R\} \rightarrow P \text{ false } R$
- token** : $P \text{ false } Token$
- \perp : $\forall \{e_1 e_2 R\} \rightarrow$
 $P e_1 R \rightarrow P e_2 R \rightarrow P (e_1 \vee e_2) R$

Actual parsers

data $P : Bool \rightarrow Set \rightarrow Set$ **where**

$_{-}?\gg_{-}$: $\forall \{e_2 R_1 R_2\} \rightarrow$
 $P \text{ true } R_1 \rightarrow (R_1 \rightarrow P e_2 R_2) \rightarrow P e_2 R_2$

$_{-}!\gg_{-}$: $\forall \{R_1 R_2\} \{e_2 : R_1 \rightarrow Bool\} \rightarrow$
 $P \text{ false } R_1 \rightarrow$
 $((x : R_1) \rightarrow \infty (P (e_2 x) R_2)) \rightarrow$
 $P \text{ false } R_2$

Tokens satisfying predicate

$sat : \forall \{R\} \rightarrow (Token \rightarrow Maybe R) \rightarrow P \text{ false } R$

$sat \{R\} p = \text{token !}\gg\gg \lambda t \rightarrow \# (ok (p t))$

where

$ok\text{-index} : Maybe R \rightarrow Bool$

$ok\text{-index } \text{nothing} = \text{false}$

$ok\text{-index } (\text{just } _) = \text{true}$

$ok : (x : Maybe R) \rightarrow P (ok\text{-index } x) R$

$ok \text{ nothing} = \text{fail}$

$ok (\text{just } x) = \text{return } x$

Derived combinators

$$\begin{aligned} _ \gg\! = _ & : \forall \{e_1 e_2 R_1 R_2\} \rightarrow \\ & P e_1 R_1 \rightarrow (R_1 \rightarrow P e_2 R_2) \rightarrow \\ & P (e_1 \wedge e_2) R_2 \end{aligned}$$

$$_ \gg\! = _ \{\text{false}\} p_1 p_2 = p_1 !\gg\! = \lambda x \rightarrow \# (p_2 x)$$

$$_ \gg\! = _ \{\text{true}\} p_1 p_2 = p_1 ?\gg\! = p_2$$

$$\begin{aligned} _ \gg _ & : \forall \{e_1 e_2 R_1 R_2\} \rightarrow \\ & P e_1 R_1 \rightarrow P e_2 R_2 \rightarrow P _ R_2 \end{aligned}$$

$$p_1 \gg p_2 = p_1 \gg\! = \lambda _ \rightarrow p_2$$

Exactly n occurrences

No corecursion:

$exactly : \forall \{e R\} (n : \mathbb{N}) \rightarrow P e R \rightarrow P ? (Vec R n)$

$exactly \text{ zero } p = \text{return } []$

$exactly (\text{suc } n) p = p \quad \ggg \lambda x \rightarrow$

$exactly n p \ggg \lambda xs \rightarrow$

$\text{return } (x :: xs)$

Exactly n occurrences

$\text{exactly-index} : \text{Bool} \rightarrow \mathbb{N} \rightarrow \text{Bool}$

$\text{exactly-index } e \text{ zero} = _$

$\text{exactly-index } e (\text{suc } n) = _$

$\text{exactly} : \forall \{e R\} (n : \mathbb{N}) \rightarrow P e R \rightarrow$
 $P (\text{exactly-index } e n) (\text{Vec } R n)$

$\text{exactly } \text{zero } p = \text{return } []$

$\text{exactly } (\text{suc } n) p = p \gg\gg \lambda x \rightarrow$
 $\text{exactly } n p \gg\gg \lambda xs \rightarrow$
 $\text{return } (x :: xs)$

Problems

Monadic interface

Index of second parser depends on result of first:

```
 $a^n b^n c^n = \text{tok "a"} \star \qquad \gg= \lambda as \rightarrow$   
  let  $n = \text{length } as$  in  
  exactly  $n$  (tok "b")  $\gg$   
  exactly  $n$  (tok "c")  $\gg$   
  return  $n$ 
```

Monadic interface

Workaround:

```
 $a^n b^n c^n$  = return zero  
| tok "a" + !>>= λ as → # (  
  let  $n = \text{length } as$  in  
  exactly  $n$  (tok "b") >>  
  exactly  $n$  (tok "c") >>  
  return  $n$ )
```

Limitations of productivity checker

- ▶ Agda requires corecursion to be guarded by *constructors*.
- ▶ Not OK:

mutual

expr = *term* *sepBy* *addOp*

term = *factor* *sepBy* *mulOp*

factor = *number*

| tok "(" >> *expr* >> tok ")"

Limitations of productivity checker

Workaround:

```
data NT : Index → Set → Set where
```

```
  expr  : NT _ _
```

```
  term  : NT _ _
```

```
  factor : NT _ _
```

```
grammar : ∀ {i R} → NT i R → P NT i R
```

```
grammar expr  = ! term  sepBy addOp
```

```
grammar term  = ! factor sepBy mulOp
```

```
grammar factor = number
```

```
      | tok "(" >> ! expr >> tok ")"
```

Limitations of productivity checker

Workaround:

grammar : $\forall \{i R\} \rightarrow NT \ i R \rightarrow P \ NT \ i R$

grammar *expr* = ! *term* *sepBy* *addOp*

grammar *term* = ! *factor* *sepBy* *mulOp*

grammar *factor* = *number*

| *tok* "(" >> ! *expr* >> *tok* ")"

! : $\forall \{e c R\} \rightarrow NT \ (e, c) R \rightarrow P \ NT \ (e, \text{step } c) R$

Limitations of productivity checker

Workaround:

$! : \forall \{e \ c \ R\} \rightarrow NT \ (e, c) \ R \rightarrow P \ NT \ (e, \text{step } c) \ R$

- ▶ Index includes *inductive* approximation of left corners (c).

- ▶ Left recursive definitions are disallowed:

grammar $nt = ! \ nt \ \Rightarrow \ c = \text{step } c$

grammar $nt = ! \ nt \ \ggg \ p$

$\Rightarrow \ c = \text{step } c \cup c'$

- ▶ Can be implemented without coinduction.
- ▶ Not quite parser *combinators*.

Conclusions

Conclusions

- ▶ Simple definition of grammars without left recursion.
- ▶ Combinator parsing with termination guarantees.
- ▶ Productivity checker needs to be improved.
- ▶ Larger examples possible:
 - ▶ Grammar scheme for mixfix operators.
- ▶ Mixed induction/coinduction is fun.

?

Bonus

“Declarative” coinductive inference systems

```
data  $_ \leq _$  :  $Ty \rightarrow Ty \rightarrow Set$  where  
  arrow :  $\infty (\tau_1 \leq \sigma_1) \rightarrow \infty (\sigma_2 \leq \tau_2) \rightarrow$   
            $\sigma_1 \rightarrow \sigma_2 \leq \tau_1 \rightarrow \tau_2$   
  refl   :  $\tau \leq \tau$   
  trans  :  $\tau_1 \leq \tau_2 \rightarrow \tau_2 \leq \tau_3 \rightarrow \tau_1 \leq \tau_3$   
  ...
```