

Total Parser Combinators

Nils Anders Danielsson (Nottingham)

ICFP, Baltimore, 2010-09-29

Total Parser Combinators *Using Mixed Induction and Coinduction*

Nils Anders Danielsson (Nottingham)

ICFP, Baltimore, 2010-09-29

Introduction

<i>expr</i>	=	<i>number</i>	
	#	<i>expr</i>	$\gg= \lambda n_1 \rightarrow$
	<i>tok</i>	'+'	$\gg= \lambda _ \rightarrow$
	<i>number</i>		$\gg= \lambda n_2 \rightarrow$
	<i>return</i>	$(n_1 + n_2)$	

Introduction

$$\begin{array}{l} \text{expr} = \# \text{ expr} \cdot \text{tok '+'} \cdot \text{number} \\ \quad | \text{ number} \end{array}$$

Introduction

expr = # *expr* · tok '+' · *number*
| *number*

Introduction

$$\begin{array}{l} \textit{expr} = \# \textit{expr} \cdot \textit{tok} \textit{'+'} \cdot \textit{number} \\ \quad | \textit{number} \end{array}$$

- ▶ Left recursive.

Introduction

$$\begin{array}{l} \textit{expr} = \# \textit{expr} \cdot \textit{tok} \textit{'+'} \cdot \textit{number} \\ \quad | \textit{number} \end{array}$$

- ▶ Left recursive.
- ▶ Parsing is guaranteed to terminate (for finite input strings).

Introduction

$$\begin{array}{l} \textit{expr} = \# \textit{expr} \cdot \textit{tok} \textit{'+'} \cdot \textit{number} \\ \quad | \textit{number} \end{array}$$

Key idea

Control the shape of parsers using a careful combination of induction and coinduction.

Introduction

$$\begin{array}{l} \textit{expr} = \# (\# \textit{expr} \cdot \# (\textit{tok} \textit{'+'})) \cdot \# \textit{number} \\ | \\ \textit{number} \end{array}$$

Key idea

Control the shape of parsers using a careful combination of induction and coinduction.

Interface (roughly)

P : *Set*

empty : P

tok : *Token* $\rightarrow P$

$-|-$: $P \rightarrow P \rightarrow P$

$- \cdot -$: $P \rightarrow P \rightarrow P$

Cyclic definitions

Want to allow cyclic definitions:

$$\begin{array}{l} \text{zeros} = \text{tok '0'} \cdot \text{zeros} \\ \quad | \text{ empty} \end{array}$$

But not all of them:

$$\text{bad} = \text{bad} \mid \text{bad}$$

Solution: Make parsers partly inductive, partly coinductive.

Mixed induction and coinduction

Inductive types

```
data List (A : Set) : Set where  
  []      : List A  
  _::__  : A → List A → List A
```

Structural recursion:

```
length : List A → ℕ  
length []      = zero  
length (x :: xs) = suc (length xs)
```

Coinductive types

data *Stream* ($A : \text{Set}$) : *Set* **where**
 $_::_ : A \rightarrow \infty (\text{Stream } A) \rightarrow \text{Stream } A$

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

$\# : A \rightarrow \infty A$
 $b : \infty A \rightarrow A$

Coinductive types

```
data Stream (A : Set) : Set where  
  _::_ : A → ∞ (Stream A) → Stream A
```

Guarded corecursion:

```
infinite : Stream ℕ  
infinite = zero :: # infinite
```

Mixed induction and coinduction

$SP\ A\ B$ represents functions of type
 $Stream\ A \rightarrow Stream\ B$:

```
data  $SP\ (A\ B : Set) : Set$  where  
  get :  $(A \rightarrow SP\ A\ B) \rightarrow SP\ A\ B$   
  put :  $B \rightarrow \infty\ (SP\ A\ B) \rightarrow SP\ A\ B$ 
```

$$SP\ A\ B \approx \nu C. \mu I. ((A \rightarrow I) + B \times C)$$

Mixed induction and coinduction

Not OK:

sink : SP A B

sink = get ($\lambda _ \rightarrow$ *sink*)

OK:

copy : SP A A

copy = get ($\lambda x \rightarrow$ put x (# *copy*))

Mixed induction and coinduction

Lexicographic guarded corecursion and
higher-order structural recursion:

$$\begin{aligned} \llbracket - \rrbracket &: SP\ A\ B \rightarrow Stream\ A \rightarrow Stream\ B \\ \llbracket \text{get } f \quad \rrbracket (a :: as) &= \llbracket f\ a \rrbracket (b\ as) \\ \llbracket \text{put } b\ sp \rrbracket as &= b :: \# (\llbracket b\ sp \rrbracket as) \end{aligned}$$

Mixed induction and coinduction

Lexicographic guarded corecursion and higher-order structural recursion:

$$\begin{aligned} \llbracket _ \rrbracket &: SP\ A\ B \rightarrow Stream\ A \rightarrow Stream\ B \\ \llbracket \text{get } f \ _ \rrbracket (a :: as) &= \llbracket f\ a \rrbracket (b\ as) \\ \llbracket \text{put } b\ sp \rrbracket as &= b :: \# (\llbracket b\ sp \rrbracket as) \end{aligned}$$

Assume that `get` were coinductive:

- ▶ *sink* would be accepted.
- ▶ $\llbracket _ \rrbracket$ would not be productive.

Back to the
parser
combinators

Choice

Hard to decide infinite choice:

$$bad = bad \mid bad$$

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

Sequencing

Problematic if p' is nullable, otherwise OK:

$$p = p \cdot p'$$

$$\frac{s_1 \in p \quad t :: s_2 \in p'}{s_1 \uparrow t :: s_2 \in p}$$

Allow the first argument to be coinductive if the second does not accept the empty string, and vice versa.

Sequencing

Problematic if p' is nullable, otherwise OK:

$$p = p \cdot p'$$

$$\frac{s_1 \in p \quad t :: s_2 \in p'}{s_1 \# t :: s_2 \in p}$$

Allow the first argument to be coinductive if the second does not accept the empty string, and vice versa.

Sequencing

Problematic if p' is nullable, otherwise OK:

$$p = p \cdot p'$$

$$\frac{s_1 \in p \quad t :: s_2 \in p'}{s_1 \uparrow t :: s_2 \in p}$$

Allow the first argument to be coinductive if the second does not accept the empty string, and vice versa.

Sequencing

Problematic if p' is nullable, otherwise OK:

$$p = p \cdot p'$$

$$\frac{s_1 \in p \quad t :: s_2 \in p'}{s_1 \uparrow t :: s_2 \in p}$$

Allow the first argument to be coinductive if the second does not accept the empty string, and vice versa.

Nullability index

Let us index parsers by their nullability:

- ▶ $P : Bool \rightarrow Set$.
- ▶ $p : P$ **true** if p accepts the empty string.
- ▶ $p : P$ **false** otherwise.

Interface

mutual

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

$empty : P \text{ true}$

$tok : Token \rightarrow P \text{ false}$

$-|_ - : P n_1 \rightarrow P n_2 \rightarrow P (n_1 \vee n_2)$

$- \cdot - : \infty \langle n_2 \rangle P n_1 \rightarrow \infty \langle n_1 \rangle P n_2 \rightarrow$
 $P (n_1 \wedge n_2)$

$\infty \langle - \rangle P : Bool \rightarrow Bool \rightarrow Set$

$\infty \langle \text{false} \rangle P n = \infty (P n)$

$\infty \langle \text{true} \rangle P n = P n$

Examples

OK:

zeros = tok '0' · # *zeros*
 | empty

Not OK:

bad = *bad* | *bad* -- Not guarded.
void = empty · *void* -- Not guarded.

Examples

OK:

$$\begin{array}{l} \text{zeros} = \# \text{ zeros} \cdot \text{tok '0'} \\ \quad \quad | \text{ empty} \end{array}$$

Not OK:

$$\begin{array}{l} \text{bad} = \text{bad} | \text{bad} \quad \text{-- Not guarded.} \\ \text{void} = \text{empty} \cdot \text{void} \quad \text{-- Not guarded.} \end{array}$$

Examples

OK:

$$\begin{array}{l} \text{zeros} = \# \text{ zeros} \cdot \text{tok '0'} \\ \quad \quad | \text{ empty} \end{array}$$

Not OK:

$$\begin{array}{ll} \text{bad} = \text{bad} | \text{bad} & \text{-- Not guarded.} \\ \text{void} = \text{empty} \cdot \# \text{void} & \text{-- Not type correct.} \end{array}$$

Example

Kleene star:

mutual

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

$p \star = \text{empty} \mid p \mid p \star$

$_ \mid : P \text{ false} \rightarrow P \text{ false}$

$p \mid = p \cdot \# (p \star)$

The argument must not accept the empty string:
the infinitely ambiguous parser $\text{empty} \star$ is
not accepted.

See the paper/code...

How is parsing implemented?

- ▶ Breadth-first algorithm:
Treat one token at a time,
compute residual recogniser using
Brzozowski derivatives.
- ▶ Combination of corecursion and recursion.
- ▶ Inefficient. Can we do better?

See the paper/code...

What about expressiveness?

- ▶ The parsers are **as expressive as possible**.
- ▶ This talk:
Every decidable language over a finite alphabet can be recognised.
- ▶ Full parser combinators:
Every finitely ambiguous decidable language can be parsed.

See the paper/code...

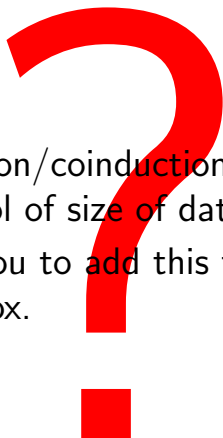
Formal semantics, algebraic laws,
mechanised correctness proofs.

Conclusions

- ▶ Mixed induction/coinduction:
Precise control of size of data.
- ▶ I encourage you to add this technique to your toolbox.

Conclusions

- ▶ Mixed induction/coinduction:
Precise control of size of data.
- ▶ I encourage you to add this technique
to your toolbox.



Bonus
slides

Infinite ambiguity

Parser combinators:

return x The empty string. Result: x .

$p \gg= f$ First p , then f applied to result of p .

fail Always fails.

▶ $(\text{return unit}) \star \mapsto [], [\text{unit}], [\text{unit}, \text{unit}], \dots$

▶ $(\text{return unit}) \star \gg= \lambda xs \rightarrow$
if $f\ xs$ **then** return unit **else** fail $\mapsto ?$

Semantics

data $_ \in _$: *List Token* $\rightarrow P\ n \rightarrow Set$ **where**

EMPTY : $[\] \in \text{empty}$

TOK : $[t] \in \text{tok } t$

ALTL : $s \in p_1 \rightarrow s \in p_1 \mid p_2$

ALTR : $s \in p_2 \rightarrow s \in p_1 \mid p_2$

SEQ : $s_1 \in b? p_1 \rightarrow s_2 \in b? p_2 \rightarrow$
 $s_1 \uparrow s_2 \in p_1 \cdot p_2$

Backend

$_ \in? _ : (s : List\ Token) (p : P\ n) \rightarrow Dec\ (s \in p)$

- ▶ Breadth-first algorithm:
Treat one token at a time,
compute residual recogniser.
- ▶ $D : (t : Token) (p : P\ n) \rightarrow P\ (D\text{-null?}\ t\ p)$.
- ▶ $t :: s \in p \Leftrightarrow s \in D\ t\ p$.
- ▶ A variant of Brzozowski's
regular expression derivatives.

Backend

$t :: s \in p \Leftrightarrow s \in D t p:$

"abc" \in $p \Leftrightarrow$

"bc" \in $D 'a' p \Leftrightarrow$

"c" \in $D 'b' (D 'a' p) \Leftrightarrow$

" " \in $D 'c' (D 'b' (D 'a' p))$

Backend

$t :: s \in p \Leftrightarrow s \in D t p:$

"abc" \in $p \Leftrightarrow$

"bc" \in $D 'a' p \Leftrightarrow$

"c" \in $D 'b' (D 'a' p) \Leftrightarrow$

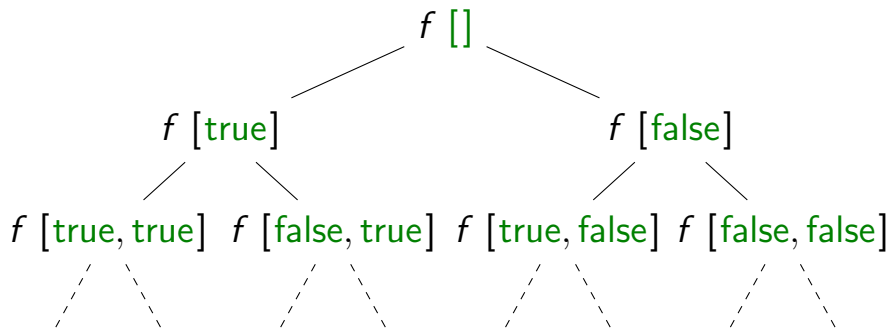
"" \in $D 'c' (D 'b' (D 'a' p))$

Can be very inefficient. Open question:
Can an efficient backend be implemented?

Expressive strength

For finite alphabets the combinators are as expressive as possible:

$f : List\ Bool \rightarrow Bool$



Expressive strength

$fail : P \text{ false}$

$fail = \# fail \cdot \# fail$

$accept\text{-if-true} : (b : Bool) \rightarrow P b$

$accept\text{-if-true true} = \text{empty}$

$accept\text{-if-true false} = fail$

$\#? : P n \rightarrow \infty \langle b \rangle P n$

$\#? \{b = \text{false}\} x = \# x$

$\#? \{b = \text{true}\} x = x$

Expressive strength

$p : (f : List\ Bool \rightarrow Bool) \rightarrow P (f [])$

$p f =$

$\# (p (\lambda xs \rightarrow f (xs ++ [true]))) \cdot \#? (tok\ true)$
 $| \# (p (\lambda xs \rightarrow f (xs ++ [false]))) \cdot \#? (tok\ false)$
 $| accept-if-true (f [])$

$s \in p f \Leftrightarrow f s \equiv true$

Laws

- ▶ The combinators form a Kleene algebra.
- ▶ Need to generalise the Kleene star:

$$\begin{aligned} _ \star &: P\ n \rightarrow P\ \text{true} \\ p \star &= (\text{nonempty } p) \star \end{aligned}$$

- ▶ The *nonempty* combinator can be defined by structural recursion:

$$\text{nonempty} : P\ n \rightarrow P\ \text{false}$$