



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]

Terminating combinator parsers in Agda

Andres Löh

based on work by Nils Anders Danielsson and Ulf Norell

Department of Information and Computing Sciences
Utrecht University

June 12, 2008

Overview

Totality

Parser combinators

Terminating combinator parsers



Totality



Total functions

A function is called **total** if it terminates and produces a valid (non- \perp) result for any input.



Total functions

A function is called **total** if it terminates and produces a valid (non- \perp) result for any input.

Many Haskell functions are not total:



Total functions

A function is called **total** if it terminates and produces a valid (non- \perp) result for any input.

Many Haskell functions are not total:

```
head :: [a] → a  
head (x : xs) = x
```

Fails on the empty list.



Total functions

A function is called **total** if it terminates and produces a valid (non- \perp) result for any input.

Many Haskell functions are not total:

```
head :: [a] → a
head (x : xs) = x
```

Fails on the empty list.

```
factorial :: Int → Int
factorial 0 = 1
factorial n = n * factorial (n - 1)
```

Loops on any negative input.



Agda is (in spirit) a total language

- ▶ All Agda functions are supposed to be total.



Agda is (in spirit) a total language

- ▶ All Agda functions are supposed to be total.
- ▶ Writing a function that the compiler cannot easily see to be terminating results in a **compiler error**.



Agda is (in spirit) a total language

- ▶ All Agda functions are supposed to be total.
- ▶ Writing a function that the compiler cannot easily see to be terminating results in a **compiler warning**.



Agda is (in spirit) a total language

- ▶ All Agda functions are supposed to be total.
- ▶ Writing a function that the compiler cannot easily see to be terminating results in a **compiler warning**.
- ▶ Other dependently typed systems (Epigram, Coq) are similar in this respect.



Agda is (in spirit) a total language

- ▶ All Agda functions are supposed to be total.
- ▶ Writing a function that the compiler cannot easily see to be terminating results in a **compiler warning**.
- ▶ Other dependently typed systems (Epigram, Coq) are similar in this respect.

Why?



Reasons for totality

- ▶ In Haskell, every type is **inhabited**: \perp is a value of any type.



Reasons for totality

- ▶ In Haskell, every type is **inhabited**: \perp is a value of any type.
- ▶ In dependently typed languages, we want to use the Curry-Howard correspondence: types are propositions, values are proofs.

data $_ \leq _ : \mathbb{N} \rightarrow \mathbb{N} \rightarrow \text{Set}$ **where**

\leq base : $\forall \{n\} \rightarrow n \leq n$

\leq step : $\forall \{m\ n\} \rightarrow m \leq n \rightarrow m \leq \text{suc } n$

trans : $\forall \{l\ m\ n\} \rightarrow l \leq m \rightarrow m \leq n \rightarrow l \leq n$

replacelnits : $\forall \{a\ m\ n\} \rightarrow m \leq n \rightarrow$
 $\text{Vec } a\ m \rightarrow \text{Vec } a\ n \rightarrow \text{Vec } a\ n$



Reasons for totality

- ▶ In Haskell, every type is **inhabited**: \perp is a value of any type.
- ▶ In dependently typed languages, we want to use the Curry-Howard correspondence: types are propositions, values are proofs.

data $_ \leq _ : \mathbb{N} \rightarrow \mathbb{N} \rightarrow \text{Set}$ **where**

\leq base : $\forall \{n\} \rightarrow n \leq n$

\leq step : $\forall \{m\ n\} \rightarrow m \leq n \rightarrow m \leq \text{suc } n$

trans : $\forall \{l\ m\ n\} \rightarrow l \leq m \rightarrow m \leq n \rightarrow l \leq n$

replacelnits : $\forall \{a\ m\ n\} \rightarrow m \leq n \rightarrow$
 $\text{Vec } a\ m \rightarrow \text{Vec } a\ n \rightarrow \text{Vec } a\ n$

- ▶ Haskell is **inconsistent**: all propositions can be proved.



Reasons for totality – contd.

- ▶ Types can contain terms:

$\text{Vec} : \text{Set} \rightarrow \mathbb{N} \rightarrow \text{Set}$

$_ \text{++ } _ : \forall \{a\ m\ n\} \rightarrow \text{Vec } a\ m \rightarrow \text{Vec } a\ n \rightarrow \text{Vec } a\ (m + n)$

$\text{tail} : \forall \{a\ n\} \rightarrow \text{Vec } a\ (\text{succ } n) \rightarrow \text{Vec } a\ n$



Reasons for totality – contd.

- ▶ Types can contain terms:

$\text{Vec} : \text{Set} \rightarrow \mathbb{N} \rightarrow \text{Set}$

$_ + _ : \forall \{a\ m\ n\} \rightarrow \text{Vec}\ a\ m \rightarrow \text{Vec}\ a\ n \rightarrow \text{Vec}\ a\ (m + n)$

$\text{tail} : \forall \{a\ n\} \rightarrow \text{Vec}\ a\ (\text{suc}\ n) \rightarrow \text{Vec}\ a\ n$

- ▶ Consider:

$\text{tail}\ (v_1 + v_2)$



Reasons for totality – contd.

- ▶ Types can contain terms:

$\text{Vec} : \text{Set} \rightarrow \mathbb{N} \rightarrow \text{Set}$

$_{-} \text{++ } _{-} : \forall \{a\ m\ n\} \rightarrow \text{Vec } a\ m \rightarrow \text{Vec } a\ n \rightarrow \text{Vec } a\ (m + n)$

$\text{tail} : \forall \{a\ n\} \rightarrow \text{Vec } a\ (\text{suc } n) \rightarrow \text{Vec } a\ n$

- ▶ Consider:

$\text{tail } (v_1 \text{++ } v_2)$

- ▶ Typechecking the expression requires unification:

$(\text{length } v_1 + \text{length } v_2) \sim \text{suc } n$

(for any n).



Consequences of totality

- ▶ Inductively defined datatypes have only **finite** values.
- ▶ Evaluation strategy (eager vs. lazy) is semantically irrelevant.
- ▶ The language cannot be Turing-complete (but still surprisingly expressive).



How to ensure totality

- ▶ Agda has a built-in **coverage** and **termination** checker.
- ▶ The coverage checker ensures that in a case analysis, all possible patterns are covered.
- ▶ The termination checker essentially checks that functions are **structurally recursive**.



Structural recursion

- ▶ Each value essentially is a constructor applied to other values:

$$| \quad v = C \ v_1 \ \dots \ v_n$$

- ▶ All such subvalues (and their subvalues ...) are **structurally smaller**. Recursive calls must make at least one argument structurally smaller.



Structural recursion

- ▶ Each value essentially is a constructor applied to other values:

$$v = C v_1 \dots v_n$$

- ▶ All such subvalues (and their subvalues ...) are **structurally smaller**. Recursive calls must make at least one argument structurally smaller.
- ▶ Many functions are trivially structurally recursive:

$$\begin{aligned} \text{length} &: \forall \{a\} \rightarrow [a] \rightarrow \mathbb{N} \\ \text{length } [] &= 0 \\ \text{length } (x :: xs) &= 1 + \text{length } xs \end{aligned}$$

Others (e.g. Quicksort) require some work ...



Parser combinators



Simple parsers

We can do better, but for this talk, we choose a naïve implementation (list of successes):

Input : Set

Input = [Char]

Parser : Set \rightarrow Set

Parser $r = \text{Input} \rightarrow [r \times \text{Input}]$



Applicative interface

$\text{fail} : \forall \{r\} \rightarrow \text{Parser } r$

$\text{fail inp} = []$

$\text{succeed} : \forall \{r\} \rightarrow r \rightarrow \text{Parser } r$

$\text{succeed } x \text{ inp} = (x, \text{inp}) :: []$



Applicative interface

$$\text{fail} : \forall \{r\} \rightarrow \text{Parser } r$$
$$\text{fail inp} = []$$
$$\text{succeed} : \forall \{r\} \rightarrow r \rightarrow \text{Parser } r$$
$$\text{succeed } x \text{ inp} = (x, \text{inp}) :: []$$
$$_ \mid _ : \forall \{r\} \rightarrow \text{Parser } r \rightarrow \text{Parser } r \rightarrow \text{Parser } r$$
$$(p \mid q) \text{ inp} = p \text{ inp} \# q \text{ inp}$$


Applicative interface – contd.

symbol : Char \rightarrow Parser Char

symbol _ [] = []

symbol x (i :: inp) = if i == x then [x, inp] else []

_ \$ _ : $\forall \{r\ s\} \rightarrow (r \rightarrow s) \rightarrow \text{Parser } r \rightarrow \text{Parser } s$

f \$ p = succeed f \star p



Applicative interface – contd.

symbol : Char \rightarrow Parser Char

symbol _ [] = []

symbol x (i :: inp) = if i == x then [x, inp] else []

_ \$ _ : $\forall \{r\ s\} \rightarrow (r \rightarrow s) \rightarrow \text{Parser } r \rightarrow \text{Parser } s$

f \$ p = succeed f \star p

- ▶ The combinators are not recursive and thus accepted as total functions by Agda.



Applicative interface – contd.

```
symbol : Char → Parser Char
symbol _ [] = []
symbol x (i :: inp) = if i == x then [x, inp] else []
_ $ _ : ∀ {r s} → (r → s) → Parser r → Parser s
f $ p = succeed f ★ p
```

- ▶ The combinators are not recursive and thus accepted as total functions by Agda.
- ▶ However, nearly all interesting grammars are cyclic, and the resulting combinator parsers recursive:

```
sum : Parser ℕ
sum = (λm _ n → m + n) $ nat ★ symbol '+' ★ sum
      | nat
```



Not all parsers terminate

nat : Parser \mathbb{N}

nat = $(\lambda n d \rightarrow n * 10 + d) \$ \text{nat} \star \text{digit}$
| digit



Not all parsers terminate

nat : Parser \mathbb{N}

nat = $(\lambda n d \rightarrow n * 10 + d) \$ \text{nat} \star \text{digit}$
| digit

many : $\forall \{a\} \rightarrow \text{Parser } a \rightarrow \text{Parser } [a]$

many p = $_ :: _ \$ p \star \text{many } p$
| succeed []

optx : Parser Char

optx = symbol 'x' | succeed ' '

optxs : Parser [Char]

optxs = many optx



The rest of this talk

- ▶ We will design parser combinators so that the resulting parsers are structurally recursive.
- ▶ Left-recursive grammars (directly and indirectly) will be type-incorrect in this library.



Terminating combinator parsers



The main idea

- ▶ Look at the following graph: nodes are parsers, an edge from one node to another indicates that a parser can directly call another (without first consuming a symbol).
- ▶ For left-recursive grammars (directly or indirectly), the graph contains cycles.
- ▶ For other grammars, the graph is cycle-free, and can be expanded into a finite tree.
- ▶ If we make this tree an index of the parser type, then left-recursive parsers are no longer type-correct.



Dependency tree

data Corners : Set **where**

leaf : Corners

node₁ : Corners → Corners

node₂ : Corners → Corners → Corners

Parser : Corners → Set → Set



Dependency tree

data Corners : Set **where**

leaf : Corners

node₁ : Corners → Corners

node₂ : Corners → Corners → Corners

Parser : Corners → Set → Set

symbol : Char → Parser leaf Char

succeed : $\forall\{r\}$ → r → Parser **leaf** r

– | – : $\forall\{r\}$ → Parser c₁ r → Parser c₂ r →
Parser (**node₂** c₁ c₂) r

– * – : $\forall\{r\ s\}$ → Parser c₁ (r → s) → Parser c₂ r →
Parser **?** s



It is important to know if a parser accepts the empty word:

Empty : Set

Empty = Bool

Parser : (Empty \times Corners) \rightarrow Set \rightarrow Set



It is important to know if a parser accepts the empty word:

Empty : Set

Empty = Bool

Parser : (Empty \times Corners) \rightarrow Set \rightarrow Set

symbol : Char \rightarrow Parser (false, leaf) Char



It is important to know if a parser accepts the empty word:

Empty : Set

Empty = Bool

Parser : (Empty \times Corners) \rightarrow Set \rightarrow Set

symbol : Char \rightarrow Parser (false, leaf) Char

succeed : $\forall \{r\} \rightarrow r \rightarrow$ Parser (true, leaf) r



It is important to know if a parser accepts the empty word:

Empty : Set

Empty = Bool

Parser : (Empty \times Corners) \rightarrow Set \rightarrow Set

symbol : Char \rightarrow Parser (false, leaf) Char

succeed : $\forall \{r\} \rightarrow r \rightarrow$ Parser (true, leaf) r

– | – : $\forall \{e_1 \ c_1 \ e_2 \ c_2 \ r\} \rightarrow$

Parser (e₁, c₁) r \rightarrow Parser (e₁, c₂) r \rightarrow

Parser (e₁ \vee e₂, node₂ c₁ c₂) r



It is important to know if a parser accepts the empty word:

Empty : Set

Empty = Bool

Parser : (Empty \times Corners) \rightarrow Set \rightarrow Set

symbol : Char \rightarrow Parser (false, leaf) Char

succeed : $\forall\{r\} \rightarrow r \rightarrow$ Parser (true, leaf) r

$- | -$: $\forall\{e_1 c_1 e_2 c_2 r\} \rightarrow$

Parser (e₁, c₁) r \rightarrow Parser (e₁, c₂) r \rightarrow

Parser (e₁ \vee e₂, node₂ c₁ c₂) r

$- * -$: $\forall\{e_1 c_1 e_2 c_2 r s\} \rightarrow$

Parser (e₁, c₁) (r \rightarrow s) \rightarrow Parser (e₁, c₂) r \rightarrow

Parser (e₁ \wedge e₂, if e₁ then node₂ c₁ c₂ else c₁) s



Not done

- ▶ What about

| Parser : (Empty \times Corners) \rightarrow Set \rightarrow Set

If, as before

| Parser $_r =$ Input $\rightarrow [r \times$ Input]

then the index information is lost!



Not done

- ▶ What about

| Parser : (Empty × Corners) → Set → Set

If, as before

| Parser _ r = Input → [r × Input]

then the index information is lost!

- ▶ We have to turn Parser into an abstract datatype:

| **data** Parser : (Empty × Corners) → Set → Set **where**

...



Not done – contd.

- ▶ Recursive definitions still pose a problem.
- ▶ Does not pass the termination checker, but still type-correct:

| $p : \text{Parser (true, leaf) Char}$

| $p = p$



Not done – contd.

- ▶ Recursive definitions still pose a problem.
- ▶ Does not pass the termination checker, but still type-correct:

$$\begin{array}{l} p : \text{Parser (true, leaf) Char} \\ p = p \end{array}$$

- ▶ Recursion must change the Corners tree!

$$\begin{array}{l} !_ : \forall\{e\ c\ r\} \rightarrow \\ \text{Parser (e, c) } r \rightarrow \text{Parser (e, node}_1\ c) r \end{array}$$

Recursion via ! fails the “occurs check”:

$$p = !p$$


Not done – contd.

- ▶ Legal cyclic definitions are still far from structurally recursive:

$$p : \text{Parser} \dots$$
$$p = \dots p \dots$$

- ▶ Turn parsers (thus Corners) into function **arguments**.
- ▶ This unfortunately has quite a few implications: we turn parser combinators and also the nonterminals of grammars into datatypes, so that we can perform case analysis in a function.



Abstract parsers

ParserType = (Empty \times Corners) \rightarrow Set \rightarrow Set₁

data Parser (nt : ParserType) : ParserType **where**

! _ : $\forall\{e\ c\ r\} \rightarrow$

nt (e, c) r \rightarrow Parser nt (e, node₁ c) r

symbol : Char \rightarrow Parser nt (false, leaf) Char

return : $\forall\{r\} \rightarrow r \rightarrow$ Parser nt (true, leaf) r

...



Grammars

Grammar : ParserType \rightarrow Set₁

Grammar nt = $\forall\{e\ c\ r\} \rightarrow nt\ (e, c)\ r \rightarrow Parser\ nt\ (e, c)\ r$

data NT : ParserType **where**

nat : NT (–, –) \mathbb{N} -- indices can be inferred!

sum : NT (–, –) \mathbb{N}

grammar : Grammar NT

grammar nat = (const 1) \$ sym '1' -- simplified

grammar sum = ($\lambda m\ _\ n \rightarrow m + n$) \$!nat * symbol '+' * !sum
| !nat



Grammars

Grammar : ParserType \rightarrow Set₁

Grammar nt = $\forall\{e\ c\ r\} \rightarrow nt\ (e, c)\ r \rightarrow Parser\ nt\ (e, c)\ r$

data NT : ParserType **where**

nat : NT (–, –) \mathbb{N} -- indices can be inferred!

sum : NT (–, –) \mathbb{N}

grammar : Grammar NT

grammar nat = (const 1) \$ sym '1' -- simplified

grammar sum = ($\lambda m\ _\ n \rightarrow m + n$) \$!nat * symbol '+' * !sum
| !nat

The definition of grammar is type correct if no left-recursion is involved. It is no longer recursive.



Interpreting the parsers

```
parse : { nt : ParserType }(g : Grammar nt)
        { e : Empty } { c : Corners } { r : Set } →
        Parser nt (e, c) r →
        LoS.Parser r -- original parser type
```

```
parse g (!p)      = parse g (g p)
```

```
parse g (symbol c) = LoS.symbol c
```

```
parse g (p1 | p2)  = LoS._ | _ (parse g p1) (parse g p2)
```

```
parse g (p1 * p2)  = LoS._ * _ (parse g p1) (parse g p2)
```



Interpreting the parsers

```
parse : { nt : ParserType } (g : Grammar nt)
       { e : Empty } { c : Corners } { r : Set } →
       Parser nt (e, c) r →
       LoS.Parser r -- original parser type
```

```
parse g (!p)      = parse g (g p)
```

```
parse g (symbol c) = LoS.symbol c
```

```
parse g (p1 | p2)  = LoS._ | _ (parse g p1) (parse g p2)
```

```
parse g (p1 * p2)  = LoS._ * _ (parse g p1) (parse g p2)
```

Is this definition structurally recursive?



Interpreting the parsers

```
parse : { nt : ParserType } (g : Grammar nt)
      { e : Empty } { c : Corners } { r : Set } →
      Parser nt (e, c) r →
      LoS.Parser r -- original parser type
```

```
parse g (!p)      = parse g (g p)
```

```
parse g (symbol c) = LoS.symbol c
```

```
parse g (p1 | p2)  = LoS._ | _ (parse g p1) (parse g p2)
```

```
parse g (p1 * p2)  = LoS._ * _ (parse g p1) (parse g p2)
```

Is this definition structurally recursive?



Interpreting the parsers

```
parse : { nt : ParserType } (g : Grammar nt)
       { e : Empty } { c : Corners } { r : Set } →
       Parser nt (e, c) r →
       LoS.Parser r -- original parser type
```

```
parse g (!p)      = parse g (g p)
```

```
parse g (symbol c) = LoS.symbol c
```

```
parse g (p1 | p2)  = LoS._ | _ (parse g p1) (parse g p2)
```

```
parse g (p1 * p2)  = LoS._ * _ (parse g p1) (parse g p2)
```

Is this definition structurally recursive?

No, in the ! case, the structure of the parser can get larger; in the * case, p₂ can have a large Corners tree.



A final refinement

We refine the Input type to keep an upper bound of the length of the input string:

Input : $\mathbb{N} \rightarrow \text{Set}$

Input n = BoundedVec Char n

Parser : $\mathbb{N} \rightarrow \mathbb{N} \rightarrow \text{Set} \rightarrow \text{Set}$

Parser m n r = Input m \rightarrow [r \times Input n]



Adapting parse

```
parse : { nt : ParserType }(g : Grammar nt)
      (n : ℕ){ e : Empty }{ c : Corners }{ r : Set } →
      Parser nt (e, c) r →
      LoS.Parser n (if e then n else pred n) r
```

...



Adapting parse

```
parse : { nt : ParserType } (g : Grammar nt)
      (n : ℕ) { e : Empty } { c : Corners } { r : Set } →
      Parser nt (e, c) r →
      LoS.Parser n (if e then n else pred n) r
```

...

```
parse g n      ( _ * _ { e1 = true } p1 p2 )
  = LoS._ * _ ( parse g n p1 ) ( parse g n p2 )
  -- ok because p1 and p2 have a smaller Corners tree
```



Adapting parse

```
parse : { nt : ParserType } (g : Grammar nt)
      (n : ℕ) { e : Empty } { c : Corners } { r : Set } →
      Parser nt (e, c) r →
      LoS.Parser n (if e then n else pred n) r
```

...

```
parse g n      ( _ * _ { e1 = true } p1 p2 )
  = LoS._ * _ ( parse g n p1 ) ( parse g n p2 )
    -- ok because p1 and p2 have a smaller Corners tree
parse g 0      ( _ * _ { e1 = false } p1 p2 )
  = LoS.fail
```



Adapting parse

```
parse : { nt : ParserType } (g : Grammar nt)
      (n : ℕ) { e : Empty } { c : Corners } { r : Set } →
      Parser nt (e, c) r →
      LoS.Parser n (if e then n else pred n) r
```

...

```
parse g n      ( _ * _ { e1 = true } p1 p2 )
  = LoS._ * _ ( parse g n p1 ) ( parse g n p2 )
  -- ok because p1 and p2 have a smaller Corners tree
```

```
parse g 0      ( _ * _ { e1 = false } p1 p2 )
  = LoS.fail
```

```
parse g (suc n) ( _ * _ { e1 = false } p1 p2 )
  = LoS._ * _ ( parse g (suc n) p1 ) ( parse↑ g n p2 )
```

```
parse↑ : ... → -- like parse, but results in ...
      LoS.Parser n n r
```



Summary

- ▶ We have shown that structurally recursive parser combinators can be implemented in Agda.
- ▶ Parsers written using this library are total. Left-recursive grammars (whether directly or indirectly) are rejected at compilation time.
- ▶ More work for the implementor, not much more work for the user, except . . .
- ▶ Defining reusable recursive derived combinators (e.g. many) requires a bit of additional trickery.
- ▶ The indices (Empty and Corners) can usually be inferred.
- ▶ Efficiency in current implementations is not too good, but in principle, not much overhead is involved – most of the indices are irrelevant at run time and can be eliminated.



Interested in Agda?

Try the seminar on
“**Dependently Typed Programming**”
(INFOMDTP)
in block 1 of 2008/2009.

