



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]

Lecture C-10: Parser Combinators - Introduction

USCS 2015

Doaitse Swierstra

Utrecht University

July 6-17, 2015

1. Domain Specific Languages



A DSL is a programming language **tailored for a particular application domain**, which captures precisely the semantics of the application domain:

- ▶ Lex / Yacc (lexing and parsing)
- ▶ L^AT_EX (for document mark-up)
- ▶ Tcl/Tk (GUI scripting)
- ▶ MatLab (numerical computations)



Advantages of using a DSL are that programs:

- ▶ are easier to understand
- ▶ are easier to write
- ▶ are easier to maintain
- ▶ can be written by non-programmers

Disadvantages:

- ▶ High start-up cost
 - ▶ design and implementation of a new language is hard
- ▶ Lack of "general purpose" features (e.g. abstraction)
- ▶ Little tool support



Embed a DSL as a library in a general purpose host language:

- ▶ inherit general purpose features from host language
 - ▶ abstraction mechanism
 - ▶ type system
- ▶ inherit compilers and tools
- ▶ good integration with host language
- ▶ many DSL's can easily be used together!



There are however also disadvantages:

- ▶ Constrained by host language
 - ▶ Syntax
 - ▶ Type system
- ▶ Implementation shines through



Haskell is a very suitable host for DSEL's:

- ▶ **Polymorphism**
- ▶ **Higher-order functions**: results can be functions describing the *denotational semantics* of an expression.
- ▶ **Lazy evaluation**:
 - ▶ allows for writing programs that analyse and transform themselves
 - ▶ gives you partial evaluation for free
- ▶ **Infix syntax** allows you to make programs look nice
- ▶ **List and monad comprehensions**
- ▶ **Type classes** for passing extra arguments in an implicit way



- ▶ Parser combinators
- ▶ Pretty printing libraries
- ▶ HaskellDB for implicitly generating SQL
- ▶ QuickCheck
- ▶ GUI libraries
- ▶ WASH/CGI – for describing HTML pages
- ▶ Haskore – for describing music
- ▶ Agent based systems
- ▶ Financial combinators – for describing financial products



Old idea:

A general purpose programming language should make it easy to write any kind of program.

New idea:

A general purpose programming language should make it possible to write very powerful libraries; even if this is a substantial effort.



2. Elementary Parser Combinators



2.1 What are parser combinators



- ▶ a collection **basic parsing functions** that recognise a piece of input
- ▶ a collection of **combinators** that build new parsers out of existing ones

For the time being parsers are just Haskell functions, but we will extend on that later.



2.2 Elementary Combinators



All libraries are based on at least the following four basic combinators:

	Haskell
alternative	$p \langle \rangle q$
composition	$p \langle * \rangle q$
terminals	<i>pSym 's'</i>
empty string	<i>pSucceed</i>



Other operators, e.g. from EBNF:

		Haskell
repetition	a^*	<i>pList a</i>
option	$a?$	'opt' a
grouping	(...)	(...)



- ▶ A parser takes a string and returns, as a proof of its success not an abstract syntax tree, but a result of type a

| $String \rightarrow a$

- ▶ A parser may not consume the whole input, so we return the unused part of the input:

| $String \rightarrow (a, String)$

- ▶ There may be many possible ways of recognizing a value of type a at the beginning of the input:

| $String \rightarrow [(a, String)]$

- ▶ Why should we only accept character strings? Everything goes:

| $[s] \rightarrow [(a, [s])]$



| `type Parser s a = [s] → [(a, [s])]`

Types

`<|>` `:: Parser s a → Parser s a → Parser s a`
`<*>` `:: Parser s a → Parser s b → Parser s (a, b)`
`<*>` `:: Parser s (b → a) → Parser s b → Parser s a`
`pSym` `:: Eq s ⇒ s → Parser s s`
`pSucceed` `:: a → Parser s a`
`pFail` `:: Parser s a`

Try to remember these types. Knowing the types is half the work when programming in Haskell.



2.3 The “List of Successes” Implementation



$$\begin{aligned}
 pFail \quad inp &= [] \\
 pSucceed \ v \ inp &= [(v, inp)] \\
 (p \langle \! \rangle q) \quad inp &= p \text{ inp } ++ q \text{ inp} \\
 pSym \ a \ (s : ss) \mid a \equiv s &= [(s, ss)] \text{ -- } !!Eq \ s \\
 pSym \ a \ _ &= [] \\
 (p \langle * \rangle q) \quad inp &= [(b2a \ b, rr) \mid (b2a, r) \leftarrow p \text{ inp} \\
 &\quad , (b, \ rr) \leftarrow q \ r \\
 &]
 \end{aligned}$$


$$p_{AB} = ((p_{Succeed} (,) \langle * \rangle p_{Sym} 'A') \langle * \rangle p_{Sym} 'B')$$

We recognize a character 'B':

$$p_{Sym} 'B'$$

Preceded by the recognition of a character 'A'

$$p_{Sym} 'A' \quad p_{Sym} 'B'$$

We now insert a “dummy” parser that returns the function (,):

$$p_{Succeed} (,) \quad p_{Sym} 'A' \quad p_{Sym} 'B'$$

Combine the result using sequential composition of parsers:

$$p_{AB} = p_{Succeed} (,) \langle * \rangle p_{Sym} 'A' \quad \langle * \rangle p_{Sym} 'B'$$



Suppose we want to deal with possibly failing computations and stay as closely as possible to the original notation; how to we deal with functions applications like $e_1 e_2$.

- ▶ both the function part e_1 and the argument part e_2 can fail to compute something
- ▶ we model this with a *Maybe*
- ▶ so we want to "apply" a *Maybe* ($b \rightarrow a$) to a *Maybe* b , and produce a *Maybe* a

func 'applyTo' *arg* = **case** *func* **of**
Just $b2a \rightarrow$ **case** *arg* **of**
Just $b \rightarrow$ *Just* ($b2a b$)
Nothing \rightarrow *Nothing*
Nothing \rightarrow *Nothing*



We capture this pattern as follows:

class *Applicative* *p* **where**

$(\langle * \rangle) :: p (b \rightarrow a) \rightarrow p b \rightarrow p a$

pure :: $a \rightarrow p a$

$(\langle \$ \rangle) :: (b \rightarrow a) \rightarrow p b \rightarrow p a$

$f \langle \$ \rangle p = \text{pure } f \langle * \rangle p$

...

instance *Applicative* *Maybe* **where**

Just $f \langle * \rangle \text{Just } v = \text{Just } (f v)$

$- \langle * \rangle - = \text{Nothing}$

pure $v = \text{Just } v$



If we now write:

$$f \langle * \rangle a_1 \langle * \rangle a_2 \langle * \rangle a_3$$

we have "overloaded" the original implicit function applications in $f a_1 a_2 a_3$.

Conclusion:

Instead applying a value of type $b \rightarrow a$ to a value of type b to result in a value of type a the operator $\langle * \rangle$ applies a p -value labelled with type $b \rightarrow a$ to a p -value labelled with type b to build a p -value labelled with type a .



Using the Haskell class system and its extensions we can denote the application also as:

$$iI f a_1 a_2 a_3 Ii$$

The symbols iI and Ii are so-called **Idiom brackets**.



Every type constructor which is in the class *Monad* can be trivially made an instance of the class *Applicative*:

instance *Monad* *m* \Rightarrow *Applicative* *m* **where**

$f \langle * \rangle v = \mathbf{do}$ *ff* $\leftarrow f$

$vv \leftarrow v$ -- no reference to *fv* in rhs

\mathbf{return} (*ff* *vv*)

pure = \mathbf{return}



The essential difference is that when using the class *Applicative* we abstain from the possibility to refer to the f -value in the second binding of the **do**-construct.

Applicative is to be preferred over *Monad*, since it allows optimisations; the second part is independent of the first part and can thus be evaluated "more statically", or even analysed independent of the run of the program!



The companion class for *Applicative* is *Alternative*:

```
class Alternative m where
  ( <|> ) :: m a -> m a -> m a
  empty :: m a

instance Alternative Maybe where
  Just l <|> _ = Just l
  _ <|> r = r
  empty = Nothing
```

Attention: For the **instance** *Alternative* (*Parser s*) the value *empty* is not the parser which recognises the empty string, but the parser that always fails!



2.4 Developing an Embedded Domain Specific Language



Because the pattern:

$$pSucceed\ f\ \langle * \rangle\ p$$

occurs so often we define

$$\langle \$ \rangle$$

$$f\ \langle \$ \rangle\ p = pSucceed\ f\ \langle * \rangle\ p$$

so we can write the previous function as:

$$pAB = (,) \langle \$ \rangle\ pSym\ 'A'\ \langle * \rangle\ pSym\ 'B'$$



Often we are not interested in parts of what we have recognized:

$semIfStat\ cond\ ifpart\ thenpart = \dots$

$pIfStat = (\lambda_c_t_e_ \rightarrow semIfStat\ c\ t\ e)$

$\langle \$ \rangle pIfToken \quad \langle * \rangle pExpr$

$\langle * \rangle pThenToken \quad \langle * \rangle pExpr$

$\langle * \rangle pElseToken \quad \langle * \rangle pExpr$

$\langle * \rangle pFiToken$

We define

$p \langle * \rangle q = (\lambda x_ \rightarrow x) \langle \$ \rangle p \langle * \rangle q$

$p \langle * \rangle q = (\lambda_ y \rightarrow y) \langle \$ \rangle p \langle * \rangle q$

$f \langle \$ \rangle q = pSucceed\ f \langle * \rangle q$

So we can now write:

$pIfStat$

$semIfStat \langle \$ \rangle pIfToken$

$\langle * \rangle pExpr$

$\langle * \rangle pThenToken \langle * \rangle pExpr$



infixl 2 *opt*

opt :: *Parser s a* → *a* → *Parser s a*

p 'opt' *v* = *p* <|> *pSucceed v*

pList :: *Parser s a* → *Parser s [a]*

pList p = (*:*) <\$> *p* <*> *pList p* 'opt' []

In the library we have special **greedy** versions which chooses the longest alternative.



Write a function that recognises a sequence of balanced parentheses, (i.e. $()$, $(())$, $((())())$, \dots , and computes the maximal nesting depth (here $1, 2, 2, \dots$). The grammar describing this language is:

$$S \rightarrow (S) S \mid \cdot$$

$$pP = (max \cdot (+1)) \langle \$ pSym \rangle (\langle * \rangle pP \langle * pSym \rangle) \langle * \rangle pP$$

'opt'

0



Parsing Expressions with Left-Associative Operators

§2.4

$pChainl :: Parser\ s\ (c \rightarrow c \rightarrow c) \rightarrow Parser\ s\ c \rightarrow Parser\ s\ c$

$pChainl\ op\ x = (f\ \langle \$ \rangle\ x\ \langle * \rangle\ pList\ (flip\ \langle \$ \rangle\ op\ \langle * \rangle\ x))$

where

$f\ x\ [] = x$

$f\ x\ (func : rest) = f\ (func\ x)\ rest$



It is not a good idea to have parsers that have alternatives starting with the same element:

$$p = \begin{array}{l} f \langle \$ \rangle q \langle * \rangle r1 \\ \langle | \rangle g \langle \$ \rangle q \langle * \rangle r2 \end{array}$$

So we define:

$$\begin{array}{l} p \langle ** \rangle q :: \text{Parser } s \ b \rightarrow \text{Parser } s \ (b \rightarrow a) \rightarrow \text{Parser } s \ a \\ p \langle ** \rangle q = (\lambda p v \ q v \rightarrow q v \ p v) \langle \$ \rangle p \langle * \rangle q \end{array}$$

So we can replace the above code by:

$$\begin{array}{l} p = q \langle ** \rangle (\text{flip } f \langle \$ \rangle r1 \langle | \rangle \text{flip } g \langle \$ \rangle r2) \\ \text{flip } f \ x \ y = f \ y \ x \end{array}$$



$\langle ?? \rangle$

$p \langle ?? \rangle q :: \text{Parser } s \ a \rightarrow \text{Parser } s \ (a \rightarrow a) \rightarrow \text{Parser } s \ a$
 $p \langle ?? \rangle q = p \langle ** \rangle (q \text{ 'opt' } id)$

$pChainr$

$pChainr \ sep \ p = p \langle ?? \rangle (flip \ \$) \ sep \ \langle * \rangle \ pChainr \ sep \ p)$

$pParens$

$pParens \quad :: \ s \rightarrow s \rightarrow \text{Parser } s \ a \rightarrow \text{Parser } s \ a$
 $pParens \ l \ r \ p = pSym \ l \ * \ p \ \langle * \ pSym \ r$



We want to recognise expressions with as result a value of the type:

```

data Expr = Lambda   Id       Expr
          | App      Expr     Expr
          | TypedExpr TypeDescr Expr
          | Ident    Id
    
```

```

pFactor = Lambda <$ pSym '\\\' <*> pIdent
           <*> pSym \.' <*> pExpr
           <|> pParens '( ' ') <*> pExpr
           <|> Ident <$> pIdent
    
```

```

pExpr = pChain1 (pSucceed App)
         (pFactor <?>) (TypedExpr <$ pTok " :: "
                        <*> pTypeDescr))
    
```



2.5 Monadic Parsers



The Chomsky hierarchy:

- ▶ Regular
- ▶ Context-free
- ▶ Context-sensitive
- ▶ Recursively enumerable

It is well known that context free grammars have limited expressibility.



times $:: \text{Int} \rightarrow \text{Parser } s \ a \rightarrow \text{Parser } s \ [a]$
 $0 \text{ 'times' } p = p\text{Succeed } []$
 $n \text{ 'times' } p = (:) \langle \$ \rangle p \langle * \rangle (n - 1) \text{ 'times' } p$
 $abc \ n = n \langle \$ \rangle (n \text{ 'times' } a)$
 $\quad \quad \quad \langle * \rangle (n \text{ 'times' } b)$
 $\quad \quad \quad \langle * \rangle (n \text{ 'times' } c)$
 $ABC = \text{foldr } (\langle \rangle) p\text{Fail } [abc \ n \mid n \leftarrow 0..]$

We admit that this is not very efficient, but left factorisation is not so easy since the corresponding context free grammar is infinite.



Recognising $\{a^n b^n c^n\}$: The Monadic Approach §2.5

Wouldn't it be nice if we could start by just recognising a sequence of a 's, and then use the result to enforce the right number of b 's and c 's?

instance *Monad* (*Parser s*) **where**

$$p \ (\gg\equiv) \ q = \lambda inp \rightarrow [(qv, rr) \mid (pv, r) \leftarrow p \quad inp \\ , (qv, rr) \leftarrow q \ pv \ r$$
$$\left. \right]$$
$$\text{return } v = \lambda inp \rightarrow [(v, inp)]$$

as :: *Parser Char Int*

as = *length* $\langle \$ \ pList \ (pSym \ 'a')$

bc n = $n \ \langle \$ \ (n \ 'times' \ b) \ \langle * \ (n \ 'times' \ c)$

ABC = **do** $n \leftarrow as$
 $bc \ n$



Can we achieve the same effect by using :

$$\{a^n b^n c^n\} = \{a^i b^j c^j\} \cap \{a^j b^j c^i\}?$$

We redefine the type *Parser* to keep track of the length of the recognized part:

$$\begin{aligned} \text{type Parser } s \ a &= [s] \rightarrow [(a, [s], \text{Int})] \\ (p \ \langle * \rangle \ q) \ \text{inp} &= [(pv \ qv, rr, pc + qc) \mid (pv, r, pc) \leftarrow p \ \text{inp} \\ &\quad , (qv, rr, qc) \leftarrow q \ r \\ &\quad] \\ &\dots \end{aligned}$$



$$\begin{aligned}
 (\langle \&\& \rangle) &:: \text{Parser } s \ a \rightarrow \text{Parser } s \ b \rightarrow \text{Parser } s \ (a, b) \\
 (p \langle \&\& \rangle q) \text{ inp} &= [((pv, qv), r, pc) \mid (pv, r, pc) \leftarrow p \text{ inp} \\
 &\quad , (qv, -, qc) \leftarrow q \text{ inp} \\
 &\quad , pc \equiv qc \\
 &]
 \end{aligned}$$

$$\begin{aligned}
 jj \ x \ y &= \text{let } xy = \quad (2+) \ \langle \$ \ pSym \ x \ \langle * \rangle \quad xy \\
 &\quad \langle * \ pSym \ y \ 'opt' \ 0 \\
 &\quad \text{in } xy \\
 i \ x &= (1+) \ \langle \$ \ pSym \ x \ \langle * \rangle \ i \ x \ 'opt' \ 0 \\
 ABC &= \quad (+) \ \langle \$ \ i \ 'a' \quad \langle * \rangle \ jj \ 'b' \ 'c' \\
 &\quad \langle \&\& \rangle \ (+) \ \langle \$ \ jj \ 'a' \ 'b' \quad \langle * \rangle \ i \quad 'c'
 \end{aligned}$$



2.6 Problems



- ▶ If your input does not conform to the language recognised by the parser all you get as a result is: `[]`.
- ▶ It may take quite a while before you get this negative result, since the backtracking may try all other alternatives at all positions.
- ▶ There is no indication of where things went wrong.

These problem have been cured in both [Parsec](#) and the [UUParsing-library](#). The latter does this:

- ▶ without much overhead
- ▶ without need for help from the programmer
- ▶ without stopping, so many errors can be found in a single run



As with any top-down parsing method, having **left-recursive parsers is no good**

- ▶ You get non-terminating parsers
- ▶ You get no error messages

This problem can be partially been cured by using chaining combinators.



Our naïve “List of successes” implementation has further drawbacks:

- ▶ The complete input has to be parsed before any result is returned
- ▶ The complete input is present in memory as long as no parse has been found
- ▶ Efficiency may depend critically on the ordering of the alternatives, and thus on how the grammar was written

For all of these problems we have found solutions in the `uu-parsinglib` package.

