

Type-safe self-inspecting code

Arthur I. Baars and S. Doaitse Swierstra

Utrecht University
{arthurb,doaitse}@cs.uu.nl

November 24, 2004



Outline

Domain-Specific-Languages

Combinator Parsers

Automatic left-recursion detection and removal

Typed abstract syntax

Constructing Grammars

Syntax Macros

Conclusions



Domain Specific Languages

A DSL is:

A programming language tailored for a particular application domain, which captures precisely the semantics of the application domain .

DSL examples:

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



Domain Specific Languages

Advantages:

- ▶ Using a DSL, programs:
 - are easier to understand
 - are quicker 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



Domain Specific *Embedded Languages*

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

- ▶ 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



DSEL's in Haskell

Haskell is a very suitable host for DSEL's:

- ▶ Polymorphism
- ▶ Lazy evaluation
- ▶ Higher-order functions
- ▶ infix syntax
- ▶ list and monad comprehension
- ▶ type classes

Examples:

- ▶ Parser combinators
- ▶ Pretty printing libraries
- ▶ HaskellDB
- ▶ QuickCheck
- ▶ GUI libraries
- ▶ WASH/CGI
- ▶ Haskore



Combinator Parsers

Parser combinators are an embedding of BNF in Haskell:

infix 6<\$>; **infix 5**<*>; **infix 4**<|>

succeed :: *a* → *Parser a*

symbol :: *Char* → *Parser Char*

<|> :: *Parser a* → *Parser a* → *Parser a*

<*> :: *Parser (a → b)* → *Parser a* → *Parser b*

failp :: *Parser a*

<\$> :: (*a* → *b*) → *Parser a* → *Parser b*

For example the BNF grammar:

$S \rightarrow 'a' S$

$| \epsilon$

is written using combinators as:

s :: *Parser [Char]*

s = (:) <\$> *symbol* 'a' <*> *s*

<|> *succeed* []



Combinator Parsers

A direct translation of:

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr } '- ' \text{ Integer} \\ &| \text{Expr } '+ ' \text{ Integer} \\ &| \text{ Integer} \end{aligned}$$

would be:

$$\begin{aligned} \text{expr} &:: \text{Parser Integer} \\ \text{expr} &= (\lambda x _ y \rightarrow x - y) \langle \$ \rangle \text{expr} \langle * \rangle \text{symbol } '- ' \langle * \rangle \text{integer} \\ &\langle | \rangle (\lambda x _ y \rightarrow x + y) \langle \$ \rangle \text{expr} \langle * \rangle \text{symbol } '+ ' \langle * \rangle \text{integer} \\ &\langle | \rangle \text{integer} \\ \text{integer} &:: \text{Parser Integer} \\ \text{integer} &= \dots \end{aligned}$$

Unfortunately, combinator parsers cannot deal with left-recursion



Combinator Parsers: *chain*

Transformed grammar, without left-recursion:

$Integer \ ((' + ' \mid ' - ') Integer) *$

Repetition combinator:

$many \ :: \ Parser \ a \ \rightarrow \ Parser \ [a]$

$many \ p \ = \ (:) \ \langle \$ \rangle \ p \ \langle * \rangle \ many \ p \ \langle | \rangle \ succeed \ []$

Chain combinator, parses as right-recursion, and converts the result to associate left

$chain \ :: \ Parser \ (a \ \rightarrow \ a \ \rightarrow \ a) \ \rightarrow \ Parser \ a \ \rightarrow \ Parser \ a$

$chain \ op \ p \ = \ f \ \langle \$ \rangle \ p \ \langle * \rangle \ many \ (flip \ \langle \$ \rangle \ op \ \langle * \rangle \ p)$

where $f \ e \ fs \ = \ foldl \ (flip \ (\$)) \ e \ fs$

Revised *expr*:

$expr \ :: \ Parser \ Integer$

$expr \ = \ chain \ op \ integer$

where $op \ = \ const \ (-) \ \langle \$ \rangle \ symbol \ ' - '$

$\ \langle | \rangle \ const \ (+) \ \langle \$ \rangle \ symbol \ ' + '$



Combinator Parsers: *chain*

The grammar can easily be extended with $*$ and $/$, by reusing the *chain* combinator:

expr :: Parser Integer

expr = chain op term

where op = const (-) <\$> symbol '-'
<|> const (+) <\$> symbol '+'

term :: Parser Integer

term = chain op integer

where op = const div <\$> symbol '/'
<|> const (*) <\$> symbol '*'



Combinator Parsers

Often left-factoring is needed to gain reasonably efficient parsers.
For example:

$$\begin{array}{l} S \rightarrow 'a' S \\ \quad | 'a' 'b' S \end{array}$$

can be rewritten to:

$$S \rightarrow 'a' (S \mid 'b' S)$$

For more complex grammars, left-factoring can be very tedious:

$$\begin{array}{l} stat \rightarrow pat ' < - ' exp \\ \quad | exp \\ \quad | ' let ' decls \end{array}$$


Combinator Parsers

Combinator parsers

- ▶ good integration with host language
- ▶ mimic BNF-syntax in Haskell using operators
- ▶ reuse through abstraction mechanism
- ▶ many cases of left-recursion can be rewritten using *chain* combinator
- ▶ complex forms of left-recursion(e.g. in mutually recursive definitions) are hard to remove
- ▶ left-factoring is often needed, and can be very tedious



Combinator Parsers

Parser Generator, such as Happy

- ▶ separate tool
- ▶ BNF like syntax
- ▶ lack abstraction mechanism, so lots of code duplication
- ▶ analyse grammar
 - perform factoring automatically
 - remove left-recursion automatically



Automatic left-recursion detection and removal

Example grammar:

```
p → q 'a'  
q → p 'b'  
   | 'c'
```

Is written using combinators as follows:

```
p, q :: Parser [Char]  
p = flip (:) <$> q <*> symbol 'a'  
q = flip (:) <$> p <*> symbol 'b'  
  <|> (:[ ]) <$> symbol 'c'
```



Automatic left-recursion detection and removal

Inlining q in p leads to:

$$p = (\lambda xs\ b\ a \rightarrow a : b : xs) \langle \$ \rangle p \langle * \rangle \text{symbol 'b'} \langle * \rangle \text{symbol 'a'} \\ \langle | \rangle (\lambda c\ a \rightarrow [c, a]) \quad \langle \$ \rangle \text{symbol 'c'} \langle * \rangle \text{symbol 'a'}$$

The parser p is left-recursive, so we need to remove the left-recursion:

$$p = f \langle \$ \rangle \text{non_left} \langle * \rangle \text{many left} \\ \textbf{where} \text{ non_left} = (\lambda c\ a \rightarrow [c, a]) \langle \$ \rangle \text{symbol 'c'} \langle * \rangle \text{symbol 'a'} \\ \text{left} = (\lambda b\ a\ xs \rightarrow a : b : xs) \langle \$ \rangle \text{symbol 'b'} \langle * \rangle \text{symbol 'a'} \\ f\ x\ fs = \text{foldl} (\text{flip } \$) x\ fs$$


Automatic left-recursion detection and removal

- ▶ Type preserving transformations
 - explicit representation of parsers: typed abstract syntax
- ▶ Make calls to other parser observable
 - Custom fix-point operator instead of normal recursion
 - Explicit representation of references



Typed abstract syntax

Datatype for representing parsers; the variable t labels each parser with its type.

```
data Parser t = Succeed t
  | Symbol Char
  | Choice (Parser t) (Parser t)
  |  $\exists x$ .Seq (Parser (x  $\rightarrow$  t)) (Parser x)
  | Fail
  |  $\exists x$ .Many (Parser x)
```

```
wrong = Succeed (+1) 'Seq' Symbol 'a'
```

The type of *Symbol 'a'* is too general, it should be *Parser Char* instead of *Parser t*.



Typed abstract syntax: smart constructors

Define smart constructors, that make sure only type-correct parser-terms can be constructed:

$succeed :: t \rightarrow Parser\ t$

$succeed = Succeed$

$symbol :: Char \rightarrow Parser\ Char$

$symbol = Symbol$

$\langle | \rangle :: Parser\ t \rightarrow Parser\ t \rightarrow Parser\ t$

$\langle | \rangle = Choice$

$\langle * \rangle :: Parser\ (a \rightarrow t) \rightarrow Parser\ a \rightarrow Parser\ t$

$\langle * \rangle = Seq$

$failp :: Parser\ t$

$failp = Fail$

$many :: Parser\ t \rightarrow Parser\ [t]$

$many = Many$

Now *wrong* is rejected:

$wrong = succeed\ (+1)\ \langle * \rangle\ symbol\ 'a'$



Compiling Parsers

Suppose we have a real parser library is a module called P

```
eval :: Parser a → P.Parser a  
eval (Succeed t) = P.succeed t  
eval (Symbol c) = P.symbol c -- wrong, why?
```



Typed abstract syntax

The datatype *Equal* is a proof that two types are equal. The only non-bottom value of this type is *self* :: *Equal a a*.

```
data Parser t =   Succeed t
                  |   Symbol (Equal Char t) Char
                  |   Choice (Parser t) (Parser t)
                  |    $\exists x.$ Seq (Parser (x  $\rightarrow$  t)) (Parser x)
                  |   Fail
                  |    $\exists x.$ Many (Equal [x] t) (Parser x)
```

symbol = *Symbol self*

many = *Many self*



Compiling Parsers, second try

Suppose we have a real parser library is a module called P

$castF :: Equal\ a\ b \rightarrow (f\ a \rightarrow f\ b)$

$eval :: Parser\ a \rightarrow P.Parser\ a$

$eval\ (Succeed\ t) = P.succeed\ t$

$eval\ (Symbol\ eq\ c) = castF\ eq\ (P.symbol\ c)$

$eval\ (Choice\ p\ q) = eval\ pP.<|>eval\ q$

$eval\ (Seq\ p\ q) = eval\ pP.<*>eval\ q$

$eval\ Fail = failp$

$eval\ Many\ eq\ p = castF\ eq\ (P.many\ p)$



Equality Type

newtype $Equal\ a\ b = Eq\ (forall\ f.f\ a \rightarrow f\ b)$

$self \quad ::\ Equal\ a\ a$

$symm \quad ::\ Equal\ a\ b \rightarrow Equal\ b\ a$

$trans \quad ::\ Equal\ a\ b \rightarrow Equal\ b\ c \rightarrow Equal\ a\ c$

$pairParts \quad ::\ Equal\ (a,b)\ (c,d) \rightarrow (Equal\ a\ c, Equal\ b\ d)$

$cast \quad ::\ Equal\ a\ b \rightarrow (a \rightarrow b)$

$castF \quad ::\ Equal\ a\ b \rightarrow (f\ a \rightarrow f\ b)$



References

```
data Ref env a
  =  $\exists$  env'.Zero (Equal env (a, env'))
  |  $\exists$ x env'.Suc (Equal env (x, env')) (Ref env' a)
```

As before we define smart constructors that pass the equality proof *self* to the corresponding constructor functions:

```
zero :: Ref (a, env) a
zero = Zero self
suc  :: Ref env' a  $\rightarrow$  Ref (b, env') a
suc  = Suc self
```

The number of *Suc*-nodes in a reference determines to which value in the environment a reference points.



References

```
deref :: Ref env a → (env → a)
deref (Zero eq)    = fst.cast eq
deref (Suc eq ref) = deref ref.snd.cast eq
```

Two arbitrary references can be compared for equality as long as they point into environments that are labeled with the same sequence of types.

```
equalRef :: Ref env a → Ref env b → Maybe (Equal a b)
equalRef (Zero eq1) (Zero eq2)
  = let (eq, _) = pairParts (inv eq1 'trans' eq2)
    in Just eq
equalRef (Suc eq1 ref1) (Suc eq2 ref2)
  = let (_, eq) = pairParts (inv eq1 'trans' eq2)
    in equalRef (cast (subF2 eq self) ref1) ref2
equalRef _ _ = Nothing
```



Grammars

A grammar is a collection of parsers. These parser may contain references to parsers in the grammar. Using nested pairs leads to infinite types.

```
grammar = (flip (:)    <$> (NT suc zero) <*> symbol 'a'  
          ,           flip (:)    <$> NT zero      <*> symbol 'b'  
            <|> (:[ ]) <$> symbol 'c'  
          )
```



Grammars

```
data Env f env
  =      EMPTY
  |  $\exists a \text{ env}' . \text{EXT } (\text{Equal } \text{env } (a, \text{env}'))$ 
      (f a) (Env f env')
```

-- smart constructors

```
empty :: Env f ()
empty = EMPTY
```

infix 1 'ext'

```
ext    :: f a  $\rightarrow$  Env f env  $\rightarrow$  Env f (a, env)
ext    = EXT self
```

```
derefEnv :: Ref env a  $\rightarrow$  (Env f env  $\rightarrow$  f a)
```

```
type Grammar env = Env (Parser env) env
```



Grammars

```
grammar :: Grammar ([Char], ([Char], ()))  
grammar = flip (:) <$> (NT suc zero) <*> symbol 'a'  
          'ext'  
          flip (:) <$> NT zero      <*> symbol 'b'  
          <|> (:[ ]) <$> symbol 'c'  
          'ext'  
empty
```



Compiling Grammars

compileGrammar :: *Grammar env* → *Env P.Parser env*

compileGrammar *gram* =

let *parsers* = *mapEnv* (*compile parsers*) *gram*

in *parsers*

mapEnv :: (*forall a.f a* → *g a*) → *Env f env* → *Env g env*

mapEnv *f* *EMPTY* = *EMPTY*

mapEnv *f* (*EXT eq x rest*) = *EXT eq (f x) (mapEnv f rest)*



Compiling Grammars

compile :: Env P.Parser env

→ Parser env a

→ P.Parser a

compile parsers prod =

case prod of

NT ref → derefEnv ref parsers

Choice p q → comp pP.<|> comp q

Symbol eq c → castF eq (P.symbol c)

Succeed x → P.succeed x

Seq p q → comp pP.<*> comp q

Many eq p → castF eq (P.many (comp p))

Fail → P.failp

where comp = compile parsers



Summary

- ▶ typed abstract syntax
- ▶ explicit references
- ▶ compile grammars to real parsers
- ▶ use of equality data type can be tedious
 - but our transformations are guaranteed to be type preserving
- ▶ loss of notational convenience when writing grammars using *zero*, and *suc*



Constructing Grammars

infixr 1 'andalso'

example :: Grammar ([Char], ([Char], ()))

example =

fixRefs

(λ~(p, (q, -)) →

flip (:) <\$> q <> symbol 'a'*

'andalso'

flip (:) <\$> p <> symbol 'b'*

<|> (:[]) <\$> symbol 'c'

'andalso'

done

)



Maybe **mdo** can help?

```
example = mdo p ← flip (:) <$> q <*> symbol 'a'  
          q ← flip (:) <$> p <*> symbol 'b' <|>  
          (:[]) <$> symbol 'c'  
          ...
```

Unfortunately, our grammars are not a *Monad*. The type of the state grows whenever we bind another parser.



Invent our own syntax

Inventing better syntax can make a combinator library much easier to use. For example Arrow syntax for the Arrows library.

```
example :: Grammar ([Char], ([Char], ()))
```

```
example = grammar
```

```
  p ← flip (:) <$> q <*> symbol 'a';
```

```
  q ← flip (:) <$> p <*> symbol 'b' <|>
```

```
    (:[ ]) <$> symbol 'c';
```

Hard-wiring special syntax into a compiler for parser combinators does not make sense.



Syntax Macros

- ▶ Generic mechanism to extend a programming language
- ▶ Defines a mapping of new concrete syntax into the core language
- ▶ Language extension modules can be loaded to facilitate a combinator library with a domain specific notation



Syntax Macros

nonterminals:

```
varid  :: String -- from Haskell Report
exp    :: Exp    -- from Haskell Report
exp10  :: Exp    -- from Haskell Report
prods  :: (Pat, Exp)
```

The new notation of grammar expressions is defined by the following macro rules:

rules:

```
exp10 ::= "grammar" (ids,ps)=prods
      => [| fixRefs (\ ~ $ids -> $ps ) |]
```

```
prods ::= => (WildP, [|done|])
```

```
prods ::= v=varid "<-" e=exp ";" (ids,ps)=prods
      => let var = VarP v
          in ( [|p ($var, $ids) |]
              , [| $e 'andalso' $ps |]
```








Conclusions

- ▶ Typed abstract syntax
 - Represent programming structures that contain internal references,
- ▶ Custom fix-point operator instead of normal recursion
 - Programs inspect their own call-graph
- ▶ Not only for combinator parsers
- ▶ Type preserving transformations
- ▶ Syntax Macros
 - provide convenient notation
 - would make Haskell an even better tool for implementing DSEL's



Related Work

-  **A. Baars and D. Swierstra.**
Syntax macros.
<http://www.cs.uu.nl/groups/ST/Center/SyntaxMacros>.
-  **A. I. Baars and S. D. Swierstra.**
Typing dynamic typing.
-  **L. Cardelli, F. Matthes, and M. Abadi.**
Extensible syntax with lexical scoping.
-  **J. Cheney and R. Hinze.**
First-Class Phantom Types.
-  **E. Pasalic and N. Linger.**
Meta-programming with typed object-language representations.

