

An Overview of *Folding Domain-Specific Languages: Deep and Shallow Embeddings*

Riley Evans (re17105)

1 Introduction

This is an overview of the techniques described in the paper *Folding Domain-Specific Languages: Deep and Shallow Embeddings*. The paper demonstrates a series of techniques that can be used when folding Domain Specific Languages. It does so through the use of a simple parallel prefix circuit language [2].

In this overview a small parser combinator language will be used. This language brings one key feature that was not described in the paper: how to apply these techniques to a typed language. Only a minimal functionally complete set of combinators have been included in the language to keep it simple. However, all other combinators usually found in a combinator language can be constructed from this set.

2 Background

2.1 DSLs

A Domain Specific Language (DSL) is a programming language that has a specialised domain or use-case. This differs from a General Purpose Language (GPL), which can be applied across a larger set of domains. DSLs can be split into two different categories: standalone and embedded. Standalone DSLs require their own compiler and typically have their own syntax. Embedded DSLs use a GPL as a host language, therefore they use the syntax and compiler from that GPL. This means that they are easier to maintain and are often quicker to develop than standalone DSLs.

An embedded DSL can be implemented with two main techniques. Firstly, a deep approach can be taken, this means that terms in the DSL will construct an Abstract Syntax Tree (AST) as a host language datatype. This can then be used to apply optimisations and then evaluated. A second approach is to define the terms as first class components of the language, avoiding the creation of an AST - this is known as a shallow embedding.

2.2 Parsers

A parser is used to convert a series of tokens into another language. For example converting a string into a Haskell datatype. Parser combinators provide a flexible approach to constructing parsers. Unlike parser generators, a combinator library is embedded within a host language: using combinators to construct the grammar. This makes it a suitable to demonstrate the techniques described in this paper for folding the DSL to create parsers.

The language is made up of 6 terms, they provide all the essential operations needed in a parser.

```
empty :: Parser a
pure  :: a      → Parser a
satisfy :: (Char → Bool) → Parser Char
try    :: Parser a      → Parser a
```

```
ap    :: Parser (a → b) → Parser a → Parser b
or    :: Parser a      → Parser a → Parser a
```

For example, a parser that can parse the characters 'a' or 'b' can be defined as,

```
aorb :: Parser Char
aorb = satisfy (≡ 'a') 'or' satisfy (≡ 'b')
```

A deep embedding of this parser language is defined in the algebraic datatype:

```
data Parser2 (a :: *) where
  Empty2 :: Parser2 a
  Pure2  :: a      → Parser2 a
  Satisfy2 :: (Char → Bool) → Parser2 Char
  Try2    :: Parser2 a      → Parser2 a
  Ap2     :: Parser2 (a → b) → Parser2 a → Parser2 b
  Or2     :: Parser2 a      → Parser2 a → Parser2 a
```

This can be interpreted by defining a function such as *size*, that finds the size of the AST used to construct the parser - this can be found in the appendix. *size* interprets the deep embedding, by folding over the datatype. See the appendix for how to add an interpretation with a shallow embedding.

3 Folds

It is possible to capture the shape of an abstract datatype through the *Functor* type class. It is possible to capture the shape of an abstract datatype as a *Functor*. The use of a *Functor* allows for the specification of where a datatype recurses. There is however one problem, a functor expressing the parser language is required to be typed. Parsers require the type of the tokens being parsed. For example a parser reading tokens that make up an expression could have the type *Parser Expr*. A functor does not retain the type of the parser, therefore it is required to define a special type class called *IFunctor*, which is able to maintain the type indices [3]. A full definition can be found in the appendix.

The shape of *Parser₂*, can be seen in *ParserF* where the *k a* marks the recursive spots.

```
data ParserF (k :: * → *) (a :: *) where
  EmptyF :: ParserF k a
  PureF  :: a      → ParserF k a
  SatisfyF :: (Char → Bool) → ParserF k Char
  TryF    :: k a      → ParserF k a
  ApF     :: k (a → b) → k a → ParserF k b
  OrF     :: k a      → k a → ParserF k a
```

```
instance IFunctor ParserF where
  imap - EmptyF    = EmptyF
  imap - (PureF x)  = PureF x
  imap - (SatisfyF c) = SatisfyF c
```

$$\begin{aligned} \text{imap } f (\text{TryF } px) &= \text{TryF } (f \text{ } px) \\ \text{imap } f (\text{ApF } pf \text{ } px) &= \text{ApF } (f \text{ } pf) (f \text{ } px) \\ \text{imap } f (\text{OrF } px \text{ } py) &= \text{OrF } (f \text{ } px) (f \text{ } py) \end{aligned}$$

Fix is used to get the fixed point of the functor. It contains the structure needed to make the datatype recursive. *Parser₄* is the fixed point of *ParserF*.

type *Parser₄* *a* = *Fix ParserF a*

A mechanism is now required for folding this abstract datatype. This is possible through the use of a catamorphism, which is a generalised way of folding an abstract datatype. Therefore, the *cata* function can be used - a definition can be found in the appendix.

Now all the building blocks have been defined that allow for the folding of the parser DSL. *size* can be defined as a fold, which is determined by the *sizeAlg*. Due to parsers being a typed language, a constant functor is required to preserve the type indicies.

```
newtype C a i = C { unConst :: a }
sizeAlg :: ParserF (C Size) a → C Size a
sizeAlg EmptyF      = C 1
sizeAlg (PureF _)   = C 1
sizeAlg (SatisfyF c) = C 1
sizeAlg (TryF (C n)) = C $ n + 1
sizeAlg (ApF (C pf) (C px)) = C $ pf + px + 1
sizeAlg (OrF (C px) (C py)) = C $ px + py + 1
size4 :: Parser4 a → Size
size4 = unConst ∘ cata sizeAlg
```

3.1 Multiple Interpretations

In DSLs it is common to want to evaluate multiple interpretations. For example, a parser may also want to know the maximum characters it will read (maximum munch). In a deep embedding this is simple, a second algebra can be defined.

```
type MM = Int
mmAlg :: ParserF (C MM) a → C MM a
mmAlg (PureF _)      = C 0
mmAlg EmptyF         = C 0
mmAlg (SatisfyF c)   = C 1
mmAlg (TryF (C px))  = C px
mmAlg (ApF (C pf) (C px)) = C $ pf + px
mmAlg (OrF (C px) (C py)) = C $ max px py
maxMunch4 :: Parser4 a → MM
maxMunch4 = unConst ∘ cata mmAlg
```

However, in a shallow embedding it is not as easy. To be able to evaluate both semantics a pair can be used, with both interpretations being evaluated simultaneously. If many semantics are required this can become cumbersome to define.

```
type Parser5 = (Size, MM)
size5 :: Parser5 → Size
size5 = fst
maxMunch5 :: Parser5 → Size
maxMunch5 = snd
smmAlg :: ParserF (C (Size, MM)) a → C (Size, MM) a
```

$$\begin{aligned} \text{smmAlg } (\text{PureF } _) &= C (1, 0) \\ \text{smmAlg } \text{EmptyF} &= C (1, 0) \\ \text{smmAlg } (\text{SatisfyF } c) &= C (1, 1) \\ \text{smmAlg } (\text{TryF } (C (s, mm))) &= C (s + 1, mm) \\ \text{smmAlg } (\text{ApF } (C (s, mm)) (C (s', mm'))) &= C (s + s' + 1, mm + mm') \\ \text{smmAlg } (\text{OrF } (C (s, mm)) (C (s', mm'))) &= C (s + s' + 1, \max mm \text{ } mm') \end{aligned}$$

Although this is an algebra, you are able to learn the shallow embedding from this, for example:

$$\begin{aligned} \text{ap}_5 \text{ } pf \text{ } px &= \text{smmAlg } (\text{ApF } pf \text{ } px) \\ \text{or}_5 \text{ } px \text{ } py &= \text{smmAlg } (\text{OrF } px \text{ } py) \end{aligned}$$

3.2 Dependent Interpretations

zygomorphisms

TODO: something in parsley. [4]

3.3 Context-sensitive Interpretations

Parsers themselves inherently require context sensitive interpretations - what you can parse will decide what you are able to parse in latter points of the parser.

Using the semantics from [5] we are able to implement a simple parser using an accumulating fold.

```
newtype Y a = Y { unYoda :: String → [(a, String)] }
```

```
yAlg :: ParserF Y a → Y a
yAlg (PureF x) = Y $ \ts. [(x, ts)]
yAlg EmptyF   = Y $ const []
yAlg (SatisfyF c) = Y $ \case
  []      . []
  (t : ts') → [(t, ts') | c t]
yAlg (TryF px) = px
yAlg (ApF (Y pf) (Y px)) = Y $ \ts.
  [(f x, ts'') | (f, ts') ← pf ts
    , (x, ts'') ← px ts']
yAlg (OrF (Y px) (Y py)) = Y $ \ts. px ts ++ py ts
parse :: Parser4 a → (String → [(a, String)])
parse = unYoda ∘ cata yAlg
```

3.4 Parameterized Interpretations

Previously we saw how to add multiple types of interpretations to a shallow embedding. We used pairs to allow us to have two interpretations. However, this doesn't extend very well to many more interpretations. Language support starts to fade for larger tuples and it will begin to become messy.

We already know that shallow embeddings are folds, so we could create a shallow embedding that is in terms of a single parameterized interpretation.

```
newtype Parser7 i = P7
  { unP7 :: ∀ a. (∀ j. ParserF a j → a j) → a i }
pure7 :: i → Parser7 i
pure7 x = P7 (λ h. h (PureF x))
empty7 :: Parser7 a
```

```

empty7 = P7 (λh. h EmptyF)
satisfy7 :: (Char → Bool) → Parser7 Char
satisfy7 c = P7 (λh. h (SatisfyF c))
try7 :: Parser7 a → Parser7 a
try7 px = P7 (λh. h (TryF (unP7 px h)))
ap7 :: Parser7 (a → b) → Parser7 a → Parser7 b
ap7 pf px = P7 (λh. h (ApF (unP7 pf h) (unP7 px h)))
or7 :: Parser7 a → Parser7 a → Parser7 a
or7 px py = P7 (λh. h (OrF (unP7 px h) (unP7 py h)))

```

```

empty3 = 1
try3 px = px + 1
ap3 pf px = pf + px + 1
or3 px py = px + py + 1
size3 :: Parser3 a → Size
size3 = id

```

```

class IFunctor f where
  imap :: (∀i. a i → b i) → f a i → f b i
newtype Fix f a = In (f (Fix f) a)
cata :: IFunctor f ⇒ (∀i. f a i → a i) → Fix f i → a i
cata alg (In x) = alg (imap (cata alg) x)

```

3.5 Implicitly Parameterized Interpretations

TODO

3.6 Modular Interpretations

TODO

References

- [1] Jeremy Gibbons and Nicolas Wu. “Folding Domain-Specific Languages: Deep and Shallow Embeddings”. In: *International Conference on Functional Programming*. Sept. 2014, pp. 339–347. DOI: 10.1145/2628136.2628138. URL: <http://www.cs.ox.ac.uk/jeremy.gibbons/publications/embedding.pdf>.
- [2] Ralf Hinze. “An Algebra of Scans”. In: *Mathematics of Program Construction*. Ed. by Dexter Kozen. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, pp. 186–210. ISBN: 978-3-540-27764-4.
- [3] Conor McBride. “Functional pearl: Kleisli arrows of outrageous fortune”. In: *Journal of Functional Programming (accepted for publication)* (2011).
- [4] Jamie Willis, Nicolas Wu, and Matthew Pickering. “Staged Selective Parser Combinators”. In: *Proc. ACM Program. Lang.* 4.ICFP (Aug. 2020). DOI: 10.1145/3409002. URL: <https://doi.org/10.1145/3409002>.
- [5] Nicolas Wu. *Yoda: A simple combinator library*. URL: <https://github.com/zenzike/yoda>.

4 Appendix

```

type Size = Int
size :: Parser2 a → Size
size Empty2 = 1
size (Pure2 _) = 1
size (Satisfy2 _) = 1
size (Try2 px) = 1 + size px
size (Ap2 pf px) = 1 + size pf + size px
size (Or2 px py) = 1 + size px + size py

```

```

type Parser3 a = Int
pure3 _ = 1
satisfy3 _ = 1

```