# 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 Languges. It does so through the use of a simple parallel prefix circuit language [3].

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 contructed 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 a 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 descriped in this paper for folding the DSL to create parsers.

The langauge 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 alegebraic 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 interpretted 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 expresing 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 indicies [4]. 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
```

```
imap f (TryF px)    = TryF (f px)
imap f (ApF pf px) = ApF (f pf) (f px)
imap f (OrF px py) = OrF (f px) (f py)
```

Fix is used to get the fixed point of the functor. It contains the structure needed to make the datatype recursive. $Parser_4$ is the fixed point of ParserF.

**type** $Parser_4$ 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.

**type** ParserAlg a i = ParserF a i $\to$ a i

**newtype** C a i = C {unConst :: a}

```
sizeAlg :: ParserAlg (C Size) i
sizeAlg EmptyF      = C 1
sizeAlg (PureF  _) = C 1
sizeAlg (SatisfyF _) = 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
```

$size_4$ :: $Parser_4$ a $\to$ Size
$size_4$ = unConst $\cdot$ 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 :: ParserAlg (C MM) i
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
```

$maxMunch_4$ :: $Parser_4$ a $\to$ MM
$maxMunch_4$ = unConst $\cdot$ 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** $Parser_5$ = (Size, MM)

$size_5$ :: $Parser_5$ $\to$ Size
$size_5$ = fst

$maxMunch_5$ :: $Parser_5$ $\to$ Size
$maxMunch_5$ = snd

```
smmAlg :: ParserAlg (C (Size, MM)) a
smmAlg (PureF _)          = C (1,     0)
smmAlg EmptyF             = C (1,     0)
smmAlg (SatisfyF c)        = C (1,     1)
smmAlg (TryF (C (s, mm))) = C (s + 1, mm)
smmAlg (ApF (C (s, mm)) (C (s', mm')))
    = C (s + s' + 1, mm + mm')
smmAlg (OrF (C (s, mm)) (C (s', mm')))
    = C (s + s' + 1, max mm mm')
```

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

$ap_5$ pf px = smmAlg (ApF pf px)
$or_5$ px py = smmAlg (OrF px py)

## 3.2 Dependent Interpretations

In a more complex parser combinator library that perform optimisations on a deep embedding, it could also be possible that there is a primary fold that depends on other secondary folds on parts of the AST. Folds such as this are named mutumorphisms [1], they can be implemented by tupling the functions in the fold. Willis et al. [5] makes use of a zygomorphism - a special case where the dependency is only one-way - to perform consumption analysis.

## 3.3 Context-sensitive Interpretations

Parsers themselves inherently require context sensitive interpretations - what can be parsed will depend on what has previously been parsed.

Using the semantics from Wu [6], an implementation can be given for a simple parser using an accumulating fold.

**newtype** Y a = Y {unYoda :: String $\to$ [(a, String)]}

```
yAlg :: ParserAlg Y i
yAlg (PureF x)   = Y $ λts → [(x, ts)]
yAlg EmptyF      = Y $ const []
yAlg (SatisfyF c) = Y $ λcase
    []       → []
    (t : ts') → [(t, ts') | c t]
yAlg (TryF px)   = px
yAlg (ApF (Y pf) (Y px)) = Y $ λts →
    [(f x, ts'') | (f, ts') ← pf ts
                 , (x, ts'') ← px ts']
yAlg (OrF (Y px) (Y py)) = Y $ λts → px ts ++ py ts
```

parse :: $Parser_4$ a $\to$ (String $\to$ [(a, String)])
parse = unYoda $\cdot$ cata yAlg

## 3.4 Parameterized Interpretations

Previously, when defining multiple interpretations in a shallow embedding, a tuple was used. However, this does not extend well when many interpretations are needed. Large tuples tend to lack good language support and will become messy to work with. It would be beneficial if a shallow embedding could have a parameter that gives it the interpretation.

Parser$_7$ allows for this approach, the shallow embedding is made up of first class functions that require an algebra argument. This algebra describes how the shallow embedding should fold the structure.

$$\textbf{newtype } \text{Parser}_7 \ i = P_7$$
$$\{\text{unP}_7 :: \forall a.(\forall j.\text{ParserF a j} \rightarrow a\ j) \rightarrow a\ i\}$$

$$\text{pure}_7 :: a \rightarrow \text{Parser}_7\ a$$
$$\text{pure}_7\ x = P_7\ (\lambda h \rightarrow h\ (\text{PureF x}))$$

$$\text{empty}_7 :: \text{Parser}_7\ a$$
$$\text{empty}_7 = P_7\ (\lambda h \rightarrow h\ \text{EmptyF})$$

$$\text{satisfy}_7 :: (\text{Char} \rightarrow \text{Bool}) \rightarrow \text{Parser}_7\ \text{Char}$$
$$\text{satisfy}_7\ c = P_7\ (\lambda h \rightarrow h\ (\text{SatisfyF c}))$$

$$\text{try}_7 :: \text{Parser}_7\ a \rightarrow \text{Parser}_7\ a$$
$$\text{try}_7\ px = P_7\ (\lambda h \rightarrow h\ (\text{TryF } (\text{unP}_7\ px\ h)))$$

$$\text{ap}_7 :: \text{Parser}_7\ (a \rightarrow b) \rightarrow \text{Parser}_7\ a \rightarrow \text{Parser}_7\ b$$
$$\text{ap}_7\ pf\ px = P_7\ (\lambda h \rightarrow h\ (\text{ApF } (\text{unP}_7\ pf\ h)\ (\text{unP}_7\ px\ h)))$$

$$\text{or}_7 :: \text{Parser}_7\ a \rightarrow \text{Parser}_7\ a \rightarrow \text{Parser}_7\ a$$
$$\text{or}_7\ px\ py = P_7\ (\lambda h \rightarrow h\ (\text{OrF } (\text{unP}_7\ px\ h)\ (\text{unP}_7\ py\ h)))$$

One benefit of this approach is that it allows the shallow embedding to be converted to a deep embedding.

$$\text{deep} :: \text{Parser}_7\ a \rightarrow \text{Parser}_4\ a$$
$$\text{deep parser} = \text{unP}_7\ \text{parser In}$$

Similarly it is possible to convert a deep embedding into a parameterised shallow embedding.

$$\text{shallow} :: \text{Parser}_4\ a \rightarrow \text{Parser}_7\ a$$
$$\text{shallow} = \text{cata shallowAlg}$$

$$\text{shallowAlg} :: \text{ParserAlg Parser}_7\ i$$
$$\text{shallowAlg } (\text{PureF x}) \quad = \text{pure}_7\ x$$
$$\text{shallowAlg EmptyF} \quad = \text{empty}_7$$
$$\text{shallowAlg } (\text{SatisfyF c}) = \text{satisfy}_7\ c$$
$$\text{shallowAlg } (\text{TryF px}) \quad = \text{try}_7\ px$$
$$\text{shallowAlg } (\text{ApF pf px}) = \text{ap}_7\ pf\ px$$
$$\text{shallowAlg } (\text{OrF px py}) = \text{or}_7\ px\ py$$

Being able to convert between both types of embedding, demonstrates that deep and parameterised shallow embeddings are inverses of each other.

## 3.5 Implicitly Parameterised Interpretations

The previous parameterised implementation still required the algebra to be specified. It would be helpful if it could be passed implicitly, if it can be determined from the type of the interpretation. This is possible in Haskell through the use of a type class.

$$\textbf{class } \text{Parser}_8\ \text{parser } \textbf{where}$$
$$\text{empty}_8 :: \text{parser a}$$
$$\text{pure}_8 \quad :: a \rightarrow \text{parser a}$$
$$\text{satisfy}_8 :: (\text{Char} \rightarrow \text{Bool}) \rightarrow \text{parser Char}$$
$$\text{try}_8 \quad :: \text{parser a} \quad \rightarrow \text{parser a}$$
$$\text{ap}_8 \quad :: \text{parser } (a \rightarrow b) \rightarrow \text{parser a} \rightarrow \text{parser b}$$
$$\text{or}_8 \quad :: \text{parser a} \quad \rightarrow \text{parser a} \rightarrow \text{parser a}$$

$$\textbf{newtype } \text{Size}_8\ i = \text{Size } \{\text{unSize} :: \text{Int}\}\ \textbf{deriving Num}$$

$$\textbf{instance } \text{Parser}_8\ \text{Size}_8\ \textbf{where}$$
$$\text{empty}_8 \quad = 1$$
$$\text{pure}_8\ \_ \quad = 1$$
$$\text{satisfy}_8\ \_ = 1$$
$$\text{try}_8\ px \quad = px + 1$$
$$\text{ap}_8\ pf\ px = \text{coerce } pf + \text{coerce } px + 1$$
$$\text{or}_8\ px\ py = px + py + 1$$

coerce allows for conversion between types that have the same runtime representation. This is the case for Size$_8$ and Int. To be able to reuse the previously defined algebras, a different type class can be defined.

$$\textbf{class } \text{Parser}_9\ \text{parser } \textbf{where}$$
$$\text{alg} :: \text{ParserAlg parser i}$$

$$\textbf{instance } \text{Parser}_9\ \text{Size}_8\ \textbf{where}$$
$$\text{alg} = \text{coerce} \cdot \text{sizeAlg} \cdot \text{imap coerce}$$

## 3.6 Modular Interpretations

There may be times when adding extra combinators would be convenient, for example adding a 'many' operator that allows for A modular technique to assembling DSLs would aid this process.

$$\textbf{data } \text{Empty}_{10}\ (k :: * \rightarrow *)\ (a :: *)\ \textbf{where}$$
$$\text{Empty}_{10} :: \text{Empty}_{10}\ k\ a$$

$$\textbf{data } \text{Pure}_{10}\ (k :: * \rightarrow *)\ (a :: *)\ \textbf{where}$$
$$\text{Pure}_{10} :: a \rightarrow \text{Pure}_{10}\ k\ a$$

$$\textbf{data } \text{Satisfy}_{10}\ (k :: * \rightarrow *)\ (a :: *)\ \textbf{where}$$
$$\text{Satisfy}_{10} :: (\text{Char} \rightarrow \text{Bool}) \rightarrow \text{Satisfy}_{10}\ k\ \text{Char}$$

$$\textbf{data } \text{Try}_{10}\ (k :: * \rightarrow *)\ (a :: *)\ \textbf{where}$$
$$\text{Try}_{10} :: k\ a \rightarrow \text{Try}_{10}\ k\ a$$

$$\textbf{data } \text{Ap}_{10}\ (k :: * \rightarrow *)\ (a :: *)\ \textbf{where}$$
$$\text{Ap}_{10} :: k\ (a \rightarrow b) \rightarrow k\ a \rightarrow \text{Ap}_{10}\ k\ b$$

$$\textbf{data } \text{Or}_{10}\ (k :: * \rightarrow *)\ (a :: *)\ \textbf{where}$$
$$\text{Or}_{10} :: k\ a \rightarrow k\ a \rightarrow \text{Or}_{10}\ k\ a$$

$$\textbf{data } (f :\!\!+\!\!: g)\ (k :: * \rightarrow *)\ (a :: *)\ \textbf{where}$$
$$L :: f\ k\ a \rightarrow (f :\!\!+\!\!: g)\ k\ a$$
$$R :: g\ k\ a \rightarrow (f :\!\!+\!\!: g)\ k\ a$$
$$\textbf{infixr } :\!\!+\!\!:$$

$$\textbf{instance } (\text{IFunctor f, IFunctor g})$$
$$\Rightarrow \text{IFunctor } (f :\!\!+\!\!: g)\ \textbf{where}$$
$$\text{imap f } (L\ x) = L\ (\text{imap f x})$$
$$\text{imap f } (R\ y) = R\ (\text{imap f y})$$

$$\textbf{type } \text{ParserF}_{10} = \text{Empty}_{10} :\!\!+\!\!: \text{Pure}_{10} :\!\!+\!\!: \text{Satisfy}_{10}$$
$$:\!\!+\!\!: \text{Try}_{10} \quad :\!\!+\!\!: \text{Ap}_{10} \quad :\!\!+\!\!: \text{Or}_{10}$$

$$\textbf{type } \text{Parser}_{10} = \text{Fix ParserF}_{10}$$

$$\text{aorb}_{10} :: \text{Parser}_{10}\ \text{Char}$$
$$\text{aorb}_{10} = \text{In } (R\ (R\ (R\ (R\ (R\ (R\ (\text{Or}_{10}$$
$$(\text{In } (R\ (R\ (L\ (\text{Satisfy}_{10}\ (\equiv\ \text{'a'})))))))$$
$$(\text{In } (R\ (R\ (L\ (\text{Satisfy}_{10}\ (\equiv\ \text{'b'})))))))))))))))$$

$$\textbf{class } (\text{IFunctor f, IFunctor g}) \Rightarrow f :\!\prec\!: g\ \textbf{where}$$
$$\text{inj} :: f\ k\ a \rightarrow g\ k\ a$$

```
instance IFunctor f ⇒ f :≺: f where
    inj = id
instance  {-# OVERLAPPING #-}  (IFunctor f, IFunctor g) ⇒ f :≺: (f :+: g) where
    inj = L
instance (IFunctor f, IFunctor g, IFunctor h, f :≺: g) ⇒ f :≺: (h :+: g) where
    inj = R · inj
```

Smart constructors:

$$empty_{10} :: (Empty_{10} :≺: f) ⇒ Fix\ f\ a$$
$$empty_{10} = In\ (inj\ Empty_{10})$$
$$pure_{10} :: (Pure_{10} :≺: f) ⇒ a → Fix\ f\ a$$
$$pure_{10}\ x = In\ (inj\ (Pure_{10}\ x))$$
$$satisfy_{10} :: (Satisfy_{10} :≺: f) ⇒ (Char → Bool) → Fix\ f\ Char$$
$$satisfy_{10}\ c = In\ (inj\ (Satisfy_{10}\ c))$$
$$try_{10} :: (Try_{10} :≺: f) ⇒ Fix\ f\ a → Fix\ f\ a$$
$$try_{10}\ px = In\ (inj\ (Try_{10}\ px))$$
$$ap_{10} :: (Ap_{10} :≺: f) ⇒ Fix\ f\ (a → b) → Fix\ f\ a → Fix\ f\ b$$
$$ap_{10}\ pf\ px = In\ (inj\ (Ap_{10}\ pf\ px))$$
$$or_{10} :: (Or_{10} :≺: f) ⇒ Fix\ f\ a → Fix\ f\ a → Fix\ f\ a$$
$$or_{10}\ px\ py = In\ (inj\ (Or_{10}\ px\ py))$$

$$aorb'_{10} :: (Or_{10} :≺: f, Satisfy_{10} :≺: f) ⇒ Fix\ f\ Char$$
$$aorb'_{10} = satisfy_{10}\ (≡ \text{'a'})\ `or_{10}`\ satisfy_{10}\ (≡ \text{'b'})$$

```
class IFunctor f ⇒ SizeAlg f where
    sizeAlg₁₀ :: f Size₈ i → Size₈ i
instance (SizeAlg f, SizeAlg g) ⇒ SizeAlg (f :+: g) where
    sizeAlg₁₀ (L x) = sizeAlg₁₀ x
    sizeAlg₁₀ (R y) = sizeAlg₁₀ y
instance SizeAlg Or₁₀ where
    sizeAlg₁₀ (Or₁₀ px py) = px + py + 1
instance SizeAlg Satisfy₁₀ where
    sizeAlg₁₀ (Satisfy₁₀ _) = 1
```

$$size_{10} :: SizeAlg\ f ⇒ Fix\ f\ a → Size_8\ a$$
$$size_{10} = cata\ sizeAlg_{10}$$
$$test :: Size$$
$$test = coerce\ (size_{10}\ (aorb'_{10} :: (Fix\ (Or_{10} :+: Satisfy_{10}))\ Char)$$

$$size_{10} :: SizeAlg\ f ⇒ Fix\ f\ a → Size_8\ a$$

# References

[1] M. Fokkinga. 1989. Tupling and Mutumorphisms.

[2] Jeremy Gibbons and Nicolas Wu. 2014. Folding Domain-Specific Languages: Deep and Shallow Embeddings. In *International Conference on Functional Programming*. 339–347. https://doi.org/10.1145/2628136.2628138

[3] Ralf Hinze. 2004. An Algebra of Scans. In *Mathematics of Program Construction*, Dexter Kozen (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 186–210.

[4] Conor McBride. 2011. Functional pearl: Kleisli arrows of outrageous fortune. *Journal of Functional Programming (accepted for publication)* (2011).

[5] Jamie Willis, Nicolas Wu, and Matthew Pickering. 2020. Staged Selective Parser Combinators. *Proc. ACM Program. Lang.* 4, ICFP, Article 120 (Aug. 2020), 30 pages. https://doi.org/10.1145/3409002

[6] Nicolas Wu. 2018. Yoda: A simple combinator library. https://github.com/zenzike/yoda

# A  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
empty₃ = 1
try₃ px = px + 1
ap₃ pf px = pf + pf + 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)
```

```
instance IFunctor Empty₁₀ where
    imap _ Empty₁₀ = Empty₁₀
instance IFunctor Pure₁₀ where
    imap _ (Pure₁₀ x) = Pure₁₀ x
instance IFunctor Satisfy₁₀ where
    imap _ (Satisfy₁₀ c) = Satisfy₁₀ c
instance IFunctor Try₁₀ where
    imap f (Try₁₀ px) = Try₁₀ $ f px
instance IFunctor Ap₁₀ where
    imap f (Ap₁₀ pf px) = Ap₁₀ (f pf) (f px)
instance IFunctor Or₁₀ where
    imap f (Or₁₀ px py) = Or₁₀ (f px) (f py)
```