



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]

Merging Parsers

IFIP 2.1, Rome meeting

Doaitse Swierstra

Department of Information and Computing Sciences
Utrecht University

Feb 7, 2012

1. History



Originally we had two libraries

The original uulib library had two modules:

1. permuting parsers
2. parsing merged lists



Parsing permuted structures

Permuting structures are abundant:

```
@inProceedings
{ BaarsLoehSwierstra2001,
  author      = { Baars, Arthur and Loeh, Andres
                  and Swierstra, S. Doaitse},
  title       = { Parsing Permutation Phrases},
  booktitle   = { Preliminary proceedings of
                  Haskell workshop 2001,
                  UU-CS-2001-23},
  year        = 2001,
  pages       = {171--182},
  editor      = {Hinze, Ralf},
}
```



Permuted structures

- ▶ The order of the elements is irrelevant
- ▶ Each item occurs exactly once
- ▶ Elements may have **different types**
- ▶ Some elements are optional



Permuted structures

- ▶ The order of the elements is irrelevant
- ▶ Each item occurs exactly once
- ▶ Elements may have **different types**
- ▶ Some elements are optional

Traditional ways of parsing such structures are clumsy.



Merged lists

Many inputs consist of a couple of merged lists, which we want to process separately:

1. Haskell: priorities, data definitions, types, classes, instances, type specifications, normal definitions
2. AG system: data definitions, attribute introductions, semantic functions, Haskell fragments



Observation

If we restrict lists to length < 1 , the parser for merged lists boils down to a permutation parser.



Observation

If we restrict lists to length < 1 , the parser for merged lists boils down to a permutation parser.

Can we generalise the way we parse merged lists to parse more general structures?



Observation

If we restrict lists to length < 1 , the parser for merged lists boils down to a permutation parser.

Can we generalise the way we parse merged lists to parse more general structures?

The aim of this talk is to present a binary combinator $\langle || \rangle$, such that $p \langle || \rangle q$ runs p and q in an interleaved way, i.e. the input is split into two sublists which are consumed by p respectively q .



2. Demo



3. Grammars



Applicative

The class Applicative describes **sequential composition** of “parsers”:

```
class Applicative p where
```

```
  (<*>) :: p (b → a) → p b → p a
```

```
  pure :: a                → p a
```

Parsers are combined using $\langle * \rangle$, where the result of the combined parser is produced by applying the result of the left operand (of type $b \rightarrow a$) to the result of the right operand (of type b).



Alternative

The class Alternative describes choice:

```
class Alternative p where  
  (<|>) :: p a → p a → p a  
  empty :: p a
```

Alternative parsers are combined using $\langle| \rangle$, and empty describes the always failing parser.



4. Grammars



Unwanted ambiguity

The following parser is ambiguous:

`pa = ...-- recognises the string "a"`

`pb = ...-- recognises the string "b"`

`p 'opt' v = p <|> pure v`

`ap = (++) <*> (pa 'opt' "x") <||> pb`

This parser will recognise "ab", "ba", "b" and "b" again, since the empty string recognisable by `(pa 'opt' "x")` can be thought to be located before or after the "b".



Unwanted ambiguity

The following parser is ambiguous:

```
pa = ...-- recognises the string "a"
```

```
pb = ...-- recognises the string "b"
```

```
p 'opt' v = p <|> pure v
```

```
ap = (++) <*> (pa 'opt' "x") <||> pb
```

This parser will recognise "ab", "ba", "b" and "b" again, since the empty string recognisable by (pa 'opt' "x") can be thought to be located before or after the "b". We decide to only include the second result.



Gram

The data type Gramm and Alt are used to represent merging parsers.

```
data Gram f a =    Gram [Alt f a] (Maybe a)
data Alt f a     =  $\forall b$ .Seq  (f (b  $\rightarrow$  a)) (Gram f b)
                  |  $\forall b$ .Bind (f b)           (b  $\rightarrow$  Gram f a)
                  |    Single (f a)
```

The first elements in the Seq, Bind and Single alternatives are parsers which are ready to be “run”, and which may not be interrupted, i.e. which accept a consecutive part of the input.



The data type Gramm and Alt are used to represent merging parsers.

```
data Gram f a =    Gram [Alt f a] (Maybe a)
data Alt f a     =  $\forall b.$ Seq  (f (b  $\rightarrow$  a)) (Gram f b)
                  |  $\forall b.$ Bind (f b)           (b  $\rightarrow$  Gram f a)
                  |    Single (f a)
```

The first elements in the Seq, Bind and Single alternatives are parsers which are ready to be “run”, and which may not be interrupted, i.e. which accept a consecutive part of the input.

- ▶ Alt f a wil not recognise the empty string
- ▶ the Maybe a part describes whether a grammar can accept the empty string, and the result if this is the case



Parsers can be lifted to Grammars

A requirement is that parsers can be split in a part recognising a non-empty string and a value to be returned when the empty string can be recognised:

```
getOneP :: Maybe (P t a) -- provided by the base parsing
getZeroP :: Maybe a      -- provided by the base parsing
mkGram :: P t a → Gram (P t) a
mkGram p = Gram (maybe [] (pure ∘ Single) (getOneP p))
              (getZeroP p)
```



5. Building Merging Parsers



Constructing parsers from Grammars

Grammars can be converted to parsers:

```
mkParserM :: (Monad f, Applicative f, ...) => Gram f a -> f a
mkParserM (Gram ls le)
  = foldr (<|>) (maybe empty pure le) (map mkParserAlt ls)
mkParserAlt (pb2a 'Seq' gb    ) = pb2a <*> mkParserM gb
mkParserAlt (pc    'Bind' c2ga) = pc >>= (mkParserM o c2ga)
mkParserAlt (Single pa      ) = pa
```





We will from now on ignore the Binds. The operator $\langle||\rangle$ follows the applicative interface:

$$\begin{aligned}
 \langle||\rangle &:: \text{Functor } f \Rightarrow \text{Gram } f (b \rightarrow a) \rightarrow \text{Gram } f b \rightarrow \text{Gram } f a \\
 &pg@(Gram pl pe) \langle||\rangle qg@(Gram ql qe) \\
 &= \text{Gram } [(\text{uncurry } \langle\$\rangle p) \text{ 'Seq' } (((,) \langle\$\rangle pp) \langle||\rangle qg) \\
 &\quad | p \text{ 'Seq' } pp \leftarrow pl] \\
 &++ [p \text{ 'Seq' } qg | \text{Single } p \leftarrow pl] \\
 &++ \text{maybe } [] (\lambda pv \rightarrow \text{map } (pv \langle\$\rangle) ql) pe \\
 &\quad \dots \text{-- similar for } qg \\
 &) (pe \langle*\rangle qe)
 \end{aligned}$$




We will from now on ignore the Binds. The operator $\langle||\rangle$ follows the applicative interface:

$$\begin{aligned}
 (\langle||\rangle) &:: \text{Functor } f \Rightarrow \text{Gram } f (b \rightarrow a) \rightarrow \text{Gram } f b \rightarrow \text{Gram } f a \\
 &pg@(Gram pl pe) \langle||\rangle qg@(Gram ql qe) \\
 &= \text{Gram } ([(\text{uncurry } \langle\$\rangle p) \text{ 'Seq' } (((,) \langle\$\rangle pp) \langle||\rangle qg) \\
 &\quad | p \text{ 'Seq' } pp \leftarrow pl] \\
 &++ [p \text{ 'Seq' } qg | \text{Single } p \leftarrow pl] \\
 &++ \text{maybe } [] (\lambda pv \rightarrow \text{map } (pv \langle\$\rangle) ql) pe \\
 &\quad \dots -- \text{similar for } qg \\
 &) \quad (pe \langle*\rangle qe)
 \end{aligned}$$

Note that this huge structure is built lazily during the actual parsing, as need arises!



6. Class Instances for Gram



Gram is a Functor

Grammars obey the conventional interface for parsers. The only difference is that they describe the break points.

instance Functor f \Rightarrow Functor (Gram f) **where**

fmap f (Gram alts e) = Gram (map (f<\$>) alts) (f <\$> e)

instance Functor f \Rightarrow Functor (Alt f) **where**

fmap a2c (pb2a 'Seq' gb) = ((a2c o) <\$> pb2a) 'Seq' gb

fmap a2c (Single pc) = Single (a2c <\$> pc)



Gram is Alternative

instance Functor $f \Rightarrow$ Alternative (Gram f) **where**
empty = Gram [] Nothing
Gram ps pe <|> Gram qs qe = Gram (ps ++ qs) (pe <|> qe)



Gram is Applicative

instance Functor f \Rightarrow Applicative (Gram f) **where**

pure a = Gram [] (Just a)

Gram l le <*> ~rg@(Gram r re)

= Gram (map ('fwdby'rg) l

++ maybe [] (\e \rightarrow map (e<\$>) r) le

) (le <*> re)

(pb2c2a 'Seq' gb) 'fwdby' gc

= (uncurry <\$> pb2c2a) 'Seq' ((,) <\$> gb <*> gc)

(Single pb2a) 'fwbby' gb

= pb2a 'Seq' gb



Conclusions

- ▶ Grammars are like parsers, but with $\langle || \rangle$ added
- ▶ Grammars are constructed lazily
- ▶ code is actually very simple
- ▶ types do the work, and tell us how to glue
- ▶ limited requirements on underlying parsing strategy



Conclusions

- ▶ Grammars are like parsers, but with $\langle||\rangle$ added
- ▶ Grammars are constructed lazily
- ▶ code is actually very simple
- ▶ types do the work, and tell us how to glue
- ▶ limited requirements on underlying parsing strategy

Constructs like:

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

look innocent, but branch infinitely! Special care is needed.

