

Preserving order in non-order preserving parsers

Rui Guerra

Arthur Baars

Doaitse Swierstra

João Saraiva

institute of information and computing sciences,
utrecht university

technical report UU-CS-2005-025

www.cs.uu.nl

Preserving order in non-order preserving parsers

Rui Guerra Arthur Baars
Doaitse Swierstra João Saraiva

Institute of Information and Computing Sciences
Utrecht University
The Netherlands
{rui,arthurb,doaitse}@cs.uu.nl

Universidade do Minho,
Departamento de Informática,
Campus de Gualtar, 4710-057 Braga, Portugal
jas@di.uminho.pt

June 21, 2005

Abstract

A *non-order preserving* (NOP) parser returns a result that cannot be mapped back to the original input, because the order of elements may have changed. For reasons of reporting errors it is desirable to preserve the initial order in some way. We show how to combine several existing parsers to construct a NOP parser that preserves the initial order by building a reordering functions while parsing. Besides other examples, we show how to use our technique in a parser for permutation phrases.

1 Introduction

There are many good reasons why the concrete syntax of a language differs from the abstract syntax. The former has been designed to make programs easy to read and write and maintain, and to make them look beautiful, whereas the latter has been designed to facilitate program analysis and code generation. Sometimes the abstract syntax closely follows the structure of the concrete syntax and it is rather straightforward to recover the original program structure from the abstract structure, but unfortunately it is more often the case that information got lost in building the abstract syntax tree. One of the detrimental effects of this kind of program transformations is that when we want to give feedback to the user about the outcome of the analysis performed on his transformed program, e.g. in the form of error messages, this is usually can no longer

be done in terms of the original program structure, but only by mentioning line numbers or similar information that was explicitly preserved for this purpose. This problem becomes more aggravating if we have an editor that is integrated with error reporting in order to give the programmer continuous feedback about the status of his program text.

A good error message should make a reference to the place in the original code where the error occurred, tells what happened, why it happened and how the error can be corrected or avoided. To be able to produce good error messages we thus need to keep somehow track of all the original information. On the other hand, we need to abstract from the language redundancy to keep the semantic analysis as simple as possible.

An example of the kind of re-orderings we want to deal with here can be found in the programming language Haskell, in which one can have type definitions, type specifications and function definitions in the program text, with no restriction on the order in which those elements appear. This has obvious advantages for the programmer since this enables him to group related definitions closely together in the program text. For semantic processing on the other hand it would be nice if the parser returned as a result three separate lists: one containing all the type definitions, a second one containing all the types for the functions to be defined, and a final one containing all the function definitions themselves.

A second example can be found in the input language for an attribute grammar system. The input basically is a long list of definitions of non-terminals, productions, attribute declarations and semantic functions, grouped in such a way that definitions that somehow belong together are lexically close. Also here semantic processing prefers a parser that returns a group of lists, one for each kind of element mentioned.

As a third example we take a small subset of the Java language, in which each **class** contains fields, methods and one constructor. These components can appear in any order. It is common after parsing to return the components in a pre-defined order. Indeed, the vast majority of parser libraries has shown that reordering results is advantageous for semantic analysis.

In this paper we show how to instrument an existing parser combinator library in such a way that information about the way a specific parser reorders elements is preserved, and may be used to map the abstract structure, that may be enriched with error messages and other feedback information, back into the original concrete structure.

The nice thing about using such combinators is that, once implemented, the techniques are available for free to everyone using the combinators: no further extra data for computing the inverse mapping has to be constructed and maintained explicitly.

We call parsers that reorder the recognized elements *non-order preserving(NOP)*. A *NOP* parser is usually defined by composing several alternatives: one for each allowed order. The reordering used in constructing the result depends on the alternative used. As the most simple example of a *NOP* parser we define a parser that recognizes two elements, of which the order does not

2 Parsing using combinator libraries

Parser combinator libraries for functional programming languages are well-known and subject of active research. Features like higher-orderness and the possibility to define new infix operators, allow parsers to be expressed in a concise and natural notation, that closely resembles the syntax of EBNF grammars. In this project we will focus on a specific kind of parsers, *NOP* parsers. For this reason we just briefly present the interface we will assume in subsequent sections. We want to stress, however, that our approach is not tied to any specific library.

We make use of a simple [4] interface that is parametrized by the result type of the parsers and assumes a list of characters as input. It can easily be implemented by straightforward list-of-successes parsers [2] and [5]. The parser interface used here is presented in figure 1; note that the precise definition of the type constructor *Parser* is not relevant, and may even differ from implementation to implementation.

The function *fail* represents the parser that always fails, whereas *succeed* never consumes any input and always returns the given argument as a result value. The parser *symbol* accepts solely its argument as input. If this character is encountered, *symbol* consumes and returns it, otherwise it fails. The \diamond operator denotes the sequential composition of two parsers in which, in order to obtain the result of the combined parser, the result of the first parser is applied to the result of the second. The operator \triangleleft expresses a choice between two parsers. Finally, the application operator \diamond is a parser transformer that can be used to apply a semantic function to a parse result. It can be defined in terms of *succeed* and \diamond . As an example of the use of these basic combinators we present the definition of the *pFoldr* combinator that will be used extensively in this presentation:

$$\begin{aligned} pFoldr & \quad \quad \quad :: (a \rightarrow b \rightarrow b) \rightarrow b \rightarrow Parser\ a \rightarrow Parser\ b \\ pFoldr\ binop\ empty\ parser & = pRest \\ & \quad \quad \quad \textbf{where } pRest = binop\ \diamond\ parser\ \diamond\ pRest\ \triangleleft\ succeed\ empty \end{aligned}$$

The function *pFoldr* is a *foldr* for parsers, i.e., it applies recursively the input *parser* until it fails, then it returns *empty*. Results are combined with the function *binop*, which has *empty* as its unit element. Many useful combinators can be built on the top of these basic ones.

3 Combining and inverting parsers

3.1 Combining and inverting two parsers

We start with a simple example. Suppose our input consists of a list of elements of type *a* and *b* in no specific order, and we want to construct a parser that returns a pair of lists: one containing all the elements of type *a* and one containing all the elements of type *b*. We start by defining a parser that recognizes

```

infixl 3 <>
infixl 4 <>, <>
symbol :: Char → Parser Char
fail    :: Parser a
succeed :: a → Parser a
(<>)     :: Parser a → Parser a → Parser a
(<>)     :: Parser (a → b) → Parser a → Parser b
(<>)     :: (a → b) → Parser a → Parser b
f <> p  = succeed f <> p
parse   :: Parser a → String → Maybe a

```

Figure 1: parser combinators

one element, either of type a or b , and that returns a function that inserts this element in the corresponding list in a pair of lists:

$$\begin{aligned}
a_or_b &:: (Parser\ a, Parser\ b) \rightarrow Parser\ ([a], [b]) \rightarrow ([a], [b]) \\
a_or_b\ (parser_a, parser_b) &= (put_a \langle \rangle parser_a \langle \rangle put_b \langle \rangle parser_b) \\
&\textbf{where } put_a\ a\ (res_a, res_b) = (a : res_a, res_b) \\
&\quad\quad\quad put_b\ b\ (res_a, res_b) = (res_a, b : res_b)
\end{aligned}$$

We can use the function $pFoldr$ to apply this parser repeatedly. When it fails or reaches the end of the input, $pFoldr$ returns a pair of empty lists $([], [])$. We define the parser combinator $pMerged$ that constructs a NOP parser from a pair of parsers.

$$\begin{aligned}
pMerged &:: (Parser\ a, Parser\ b) \rightarrow Parser\ ([a], [b]) \\
pMerged\ (parser_a, parser_b) &= pFoldr\ apply\ ([], [])\ (a_or_b\ (parser_a, parser_b)) \\
apply\ element\ rest &= element\ rest
\end{aligned}$$

As we can see, the put_a function inserts the value returned by $parser_a$ in the list res_a , that contains all the values of type a . The put_b function does the same as put_a but for values of type b . Suppose now that $parser_{int}$ and $parser_{alp}$ are parsers that recognize and return digits and letters, respectively. Then the NOP parser that combines these parsers is represented by:

$$\begin{aligned}
int_alp &:: Parser\ ([Int], [Char]) \\
int_alp &= pMerged\ (parser_{int}, parser_{alp})
\end{aligned}$$

The parser int_alp recognizes a list of integers and characters and returns them in a pair of lists. Thus the result of applying int_alp to the input string "A1bC2" is the pair $([1, 2], "AbC")$. Using the parser combinator $pMerged$ we can define any NOP made up from two simple parsers.

Having such a parser, the question that arises now is how to recover the original order of the input? The obvious answer is to build a function that memorizes

that order and it is able to map back the original input. In other words, we are looking for an *invert* function that takes an input like $([1, 2], \text{"AbC"})$ and reconstructs the initial order "A1bC2" .

The function that memorizes the order of the original input has to be constructed when the result is re-ordered. That is to say, it has to be constructed during parsing. So the parser function has to perform three tasks:

- Detecting whether the input follows the syntactic rules specified by underlying grammar, or not;
- Giving semantic meaning to the input (which can preserve the input order or not);
- Constructing an inverse function that recovers the initial order of the input (only in the case of a *NOP* parser).

To recover the initial order of the input, we have to construct an *invert* function, during parsing. The *invert* function will take a pair of lists and reorders its elements into a single list. Standard lists, however, require that all elements have the same type. This is not the case in our running example. To simplify our introductory explanation we take for granted that the two lists have the same type. We will return to this subject in section 5. The *invert* function takes a pair of lists and returns them in the initial order and we use a type synonym to represent its type:

type *Invert* $s = ([s], [s]) \rightarrow [s]$

The *invert* function has to remember in which side of the pair the parsing result was stored. When the semantic function inserts the parser result in the left side of the pair, the *invert* function will take an element from the left side and otherwise it will take an element from the right side. We define two *put* and two *get* functions. As before, we combine each parser with a *put* function using the combinator $\langle \diamond \rangle$. Each *put* function will update the *inv* function using the *get* function. The functions get_a and get_b update the *invert* function by taking an element from the left side of the input or from the right side, respectively. For a parser that failed or reached the end of the input, the resulting *invert* function returns an empty list.

$pMerged :: (Parser\ a, Parser\ b) \rightarrow Parser\ ([a], [b], Invert\ s)$
 $pMerged\ (parser_a, parser_b) = pFoldr\ (\$)\ ([], [])\ const\ []$
 $(put_a \langle \diamond \rangle parser_a \langle \diamond \rangle put_b \langle \diamond \rangle parser_b)$

where

$put_a\ a\ ((res_a, res_b), inv) = ((a : res_a, res_b), get_a\ inv)$
 $put_b\ b\ ((res_a, res_b), inv) = ((res_a, b : res_b), get_b\ inv)$

$get_a\ inv\ ([], res_b) = inv\ ([], res_b)$
 $get_a\ inv\ (a : res_a, res_b) = a : inv\ (res_a, res_b)$
 $get_b\ inv\ (res_a, []) = inv\ (res_a, [])$
 $get_b\ inv\ (res_a, b : res_b) = b : inv\ (res_a, res_b)$

The resulting *NOP* parser returns a pair where the first element is another pair with the recognized symbols and the second element is the *invert* function.

Suppose that we want to build a *NOP* parser that recognizes letters and distinguishes upper case from lower case letters. If *parser_{low}* and *parser_{upp}* are parsers that recognize lower case and upper case letters, respectively, then our parser is easily represented by:

$$\begin{aligned} low_upp &:: Parser ((String, String), Invert a) \\ low_upp &= pMerged (parser_{low}, parser_{upp}) \end{aligned}$$

The parser *low_upp* besides returning the two lists with the parser results, it also returns the *invert* function. Thus applying the parser *low_upp* to the string "AaBb" produces the pair (("ab", "AB"), <Invert>).

To recover the initial input, we only have to apply the invert function to the pair of recognized values.

We should realize that the *invert* function remembers the order of the symbols but not the symbols themselves. As a result, we can change the initial values and still recover the initial order. Thus if we change the value 'B' in the previous result pair for a 'Z', for instance, the inverted input will have a 'Z' in the place of the 'B'.

Note that abstractly, this can be seen as a typical error correcting task, where the parser replaces a symbol, occurring in the wrong position, by the expected one.

3.2 Combining *n* parsers

In the previous section, we saw how to combine and invert two parsers. However, a *NOP* parser can be build up with a combination of any number of parsers. In this section, we define a parser combinator that may be used to combine any number of parsers into a *nop* parser.

When we combine two parsers, for instance *Parser a* and *Parser b*, the parsing result of the new parser is a pair of type (a, b) . If we combine this new parser with another parser, for instance, *Parser c*, the parsing result, following the same logic, is of type (a, b, c) . The problem of this approach is that the type of the resulting parser depends on the number of combined parsers. Our aim is to define a function that combines a variable number of parsers. Because the types of the combined parsers are typically not the same, we cannot represent the result by lists. A nested cartesian product is more suitable for this situation. We use combinators only in a left associative way, so all products are of the type $((x, y), z)$ instead of (x, y, z) ¹.

In the previous section, both parser results were returned as lists. This is, however, a severe limitation. It would be convenient to abstract from the result of the parsers, in order to allow the use of data structures other than lists.

¹This was also the very first approach taken in the design of parser combinator libraries [2]

In constructing lists we have used the predefined insert function on lists, and the (*const []*) to construct the empty list. Thus in order to use a custom data structure we have to provide its constructors, i.e., the function that constructs the empty data structure and the respective insert function. We provide this information in a pair like (*empty, insert* \diamond *parser*). Now the parser result is a function that knows how to insert the recognized value in the data structure. The information pair can be represented by a type synonym:

type *ParseInf* *res* = (*res*, *Parser* (*res* \rightarrow *res*))

To return the parsing result in the predefined type *lists*, we have to build the information pair with the insert function (*:*) and with the empty structure (*[]*):

listOf $\quad\quad\quad$ **::** *Parser* *res* \rightarrow *ParseInf* [*res*]
listOf parser = (*[]*, (*:*) \diamond *parser*)

Any data structure can be used. Thus, for instance, we can define a data type for balanced binary trees.

data *Tree* *a* = *Leaf* | *Node* *Int* (*Tree* *a*) *a* (*Tree* *a*)

Now we only need to compose the information pair, where the *Leaf* is the empty structure and the *insert* function is the balanced insert function for binary trees.

treeOf $\quad\quad\quad$ **::** *Parser* *res* \rightarrow *ParseInf* (*Tree* *res*)
treeOf parser = (*Leaf*, *insert* \diamond *parser*)

Now instead of combining parsers, we have to combine information pairs that contain parsers. To be able to construct one parser that combines all the other parsers, we go through a two phase process. First we perform an abstract interpretation in which we combine all the information pairs into a single structure. Second we construct the parser we are interested in, using the result pair of the first phase.

To accomplish the first phase, we define a new combinator \diamond , that combines the information pairs out of which the *NOP* parser is to be constructed. First we combine the empty structures in a pair. The result of applying the parser combinator \diamond to a variable number of pairs, is a nested pair of empty structures. Next we have to compose the parser results. Following the approach in which we combined a parser with an insert function, the parser result becomes a function that knows how to insert the recognized value in the respective data structure. Thus we only apply the parsing result to the correct side of the pair.

(\diamond) $\quad\quad\quad$ **::** *ParseInf* *a* \rightarrow *ParseInf* *b* \rightarrow *ParseInf* (*a*, *b*)
(*e_a*, *p_a*) \diamond (*e_b*, *p_b*) = ((*e_a*, *e_b*), *comb_a* \diamond *p_a* \diamond *comb_b* \diamond *p_b*)
where
comb_a *insert_a* (*res_a*, *res_b*) = (*insert_a* *res_a*, *res_b*)
comb_b *insert_b* (*res_a*, *res_b*) = (*res_a*, *insert_b* *res_b*)

The result of combining all parsers is a single information pair, of which the first element is a nested pair with all the empty structures and the second element is

and in $comb_b$ we apply it to $values_b$. This function returns a pair of which the first element is the data structure updated with the new parsed element and the second one the get function that removes the corresponding element from the data structure. Secondly we update the $invert$ function by applying the get function to the correct side of the pair. The result is a pair with the recovered element and the non inverted elements represented by $(e, (rest, inp_b))$, in the case of $comb_a$ or $(e, (inp_a, rest))$ in the case of $comb_b$.

$$\begin{aligned}
comb_a \text{ insert } (values_a, values_b) &= ((newVal_a, values_b), invert_a) \\
\text{where} & \\
(newVal_a, get_a) &= \text{insert } values_a \\
invert_a (inp_a, inp_b) &= \mathbf{let} (e, rest) = get_a \text{ } inp_a \\
&\quad \mathbf{in} (e, (rest, inp_b)) \\
\\
comb_b \text{ insert } (values_a, values_b) &= ((values_a, newVal_b), invert_b) \\
\text{where} & \\
(newVal_b, get_b) &= \text{insert } values_b \\
invert_b (inp_a, inp_b) &= \mathbf{let} (e, rest) = get_b \text{ } inp_b \\
&\quad \mathbf{in} (e, (inp_a, rest))
\end{aligned}$$

In the previous section, in order to obtain the combined parser, we only needed to use the parser combinator $pFoldr$ with the function $apply$. This was possible because each parser result was a function that knew how to insert its result in the data structure.

Now we need to apply the result function to the data structure and update the invert function, as well. The $acceptNextVal$ function applies the result function $insRes$ to the data structure $recognVals$. The result is a pair with the updated data structure $newValues$ and the function $getVal$. This function is used to get the first inverted value and subsequently we append it to the rest of the inverted values computed by $getRest$.

$$\begin{aligned}
acceptNextVal &:: (a \rightarrow (a, c \rightarrow (d, c))) \rightarrow (a, c \rightarrow [d]) \rightarrow (a, c \rightarrow [d]) \\
acceptNextVal \text{ insert } (recognVals, getRest) &= (newValues, invert) \\
\text{where} & \\
(newValues, getVal) &= \text{insert } recognVals \\
invert \text{ inp} &= \mathbf{let} (v, rest) = getVal \text{ } inp \\
&\quad \mathbf{in} v : (getRest \text{ } rest)
\end{aligned}$$

The function $pFoldr$ applies $acceptNextVal$ to the result of the combined parsers, $alternatives$. When the parser fails or reaches the end of the input, $pFoldr$ returns a pair with the empty structure and a function that always returns an empty list, the invert function.

$$\begin{aligned}
pMerged &:: ParseInf a c d \rightarrow Parser (a, c \rightarrow [d]) \\
pMerged (units, alter) &= pFoldr acceptNextVal (units, const []) alter
\end{aligned}$$

The parser dig_low_upp is defined in a similar way as in the previous section,

but now we get the *invert* function, as well.

$$\begin{aligned} \text{dig_low_upp} &:: \text{Parser } ((([\text{Char}], [\text{Char}]), \text{Tree Char}), \\ &\quad ([a], [a]), \text{Tree } a) \rightarrow [a]) \\ \text{dig_low_upp} &= \text{pMerged } (\text{listOf parser}_{\text{dig}} \diamond \\ &\quad \text{listOf parser}_{\text{low}} \diamond \\ &\quad \text{treeOf parser}_{\text{upp}}) \end{aligned}$$

The usage of the parser is the same. Its result is a pair where the first element is a nested pair of parsing results and the second element the *invert* function. Thus applying the parser *dig_low_upp* to the input string "A1bC2" results in $((("12", "b"), \text{Node } 2 (\text{Node } 1 \text{ Leaf 'A' Leaf}) 'C' \text{ Leaf}), <\text{invert}>)$.

The *invert* function is able to recover the initial order from the parsing result, by just applying it to the nested pair of recognized values. Thus, to recover the input string, we simply apply the *invert* function to the nested pair $((("12", "b"), \text{Node } 2 (\text{Node } 1 \text{ Leaf 'A' Leaf}) 'C' \text{ Leaf})$.

4 Customizing the parser result

In the previous sections, we implemented a function, that combines a variable number of parsers. To achieve our aim we constructed the parsing result as a nested pair. Consequently the input of the *invert* function should be a nested pair, as well. This data type is useful to return a variable number of results, but it is difficult to understand by human beings. It would be convenient if the user could specify the type of the parsing result and the *invert* input. Therefore, we will now develop a variant in which the nested pairs are only used internally and thus eliminate the inconvenience of showing them to the outside world.

To be able to customize the parsing result, first we have to 'unpack' the nested pair, i.e., to take each element out of the nested pair. Next we have to transform the result to the custom one. Taking as example the previous combination of three parsers, we realize that the parsing result is a nested pair of three elements. To customize it we need a function that takes each different result and returns it in the new type. If for instance, we want to return this parsing result as a tuple, we have to supply the $(, ,)^2$ function. This function takes three separate arguments and returns them in a tuple. In the approach followed thus far we actually need a function of the type $((a, b), c) \rightarrow (a, b, c)$. Haskell has defined the *uncurry*³ function. It takes a function with two arguments and returns a function that takes a pair of arguments instead. If we apply the *uncurry* function to $(, ,)$, we get:

$$\text{uncurry } (, ,) :: (a, b) \rightarrow c \rightarrow (a, b, c)$$

Composing two *uncurry* function we construct the function we need:

$$(\text{uncurry}.\text{uncurry}) (, ,) :: ((a, b), c) \rightarrow (a, b, c)$$

² $(, ,) :: a \rightarrow b \rightarrow c \rightarrow (a, b, c)$

³ $\text{uncurry} :: (a \rightarrow b \rightarrow c) \rightarrow (a, b) \rightarrow c$

By applying two *uncurry* functions to the $(,)$ function we manage to adapt it to handle nested pairs. The input of the *invert* function is a nested pair as well, $invert :: ((l, l), l) \rightarrow [l]$. We need a *pack* function that adapts it to take three inputs instead of a nested pair. The type of the invert function becomes $invert :: l \rightarrow l \rightarrow l \rightarrow [l]$. Haskell has defined the *curry*⁴ function, as well. It takes a function with one pair of arguments and returns another function that takes two arguments. Like in the previous case if we compose two *curry* function, we are able to customize the invert function: The *curry* function transforms a function, that takes one pair of arguments, to another one that accepts flattened inputs.

$$(curry.curry) invert :: l \rightarrow l \rightarrow l \rightarrow [l]$$

The definition of both functions, *pack* and *unpack*, depends on the number of parsers we combine. Thus they have to be constructed during the merge process, and the *listOf* and the *treeOf* functions have to include a *pack* and an *unpack* function too. Since they store a single parser, the result does not need to be packed or unpacked, thus the *unpack* and *pack* function are defined as the *identity* function.

$$\begin{aligned} listOf\ parser &= ([], (\lambda x\ l \rightarrow (x : l, getVal)) \diamond parser, id, id) \\ treeOf\ parser &= (Leaf, insert \diamond parser, id, id) \end{aligned}$$

Whenever we combine two parsers, we have to compose the *pack* function with another *curry* function. We do the same thing to the *unpack* function, but instead of using the *curry* function we use the *uncurry*. Initially both *pack* and *unpack* functions are the *identity* function. After combining the two first parsers, the resulting *pack* and *unpack* are equal to *curry* and *uncurry* functions respectively, as we saw before. In the end we will have $(n - 1)$ *curry* and *uncurry* composed functions, where n is the number of parsers combined.

$$\begin{aligned} (e_a, p_a, unpack_a, pack_a) \diamond (e_b, p_b, unpack_b, pack_b) = & \\ ((e_a, e_b) & \\ , comb_a \diamond p_a \diamond comb_b \diamond p_b & \\ , unpack_b.uncurry.unpack_a & \\ , pack_b.curry.pack_a & \\) & \end{aligned}$$

Now we eliminate the nested pair of the input of the *invert* function by applying the *pack* function to it. The semantic function *sem* that returns the parsing result in a predefined data type is adapted to the nested pair by applying the *unpack* function to it.

$$\begin{aligned} sem\ 'pMerged'\ (units, alternatives, unpack, pack) = & \\ (\lambda(result, invert) \rightarrow (unpack\ sem\ result, pack\ invert)) & \\ \diamond & \\ pFoldr\ acceptNextVal\ (units, const\ [])\ alternatives & \end{aligned}$$

⁴*curry* :: ((a, b) → c) → a → b → c

If for instance, we want to change the previous example, to return the parser result in a tuple, instead of a nested pair, the definition of *dig_low_upp* will be:

$$dig_low_upp = (, ,) \text{ 'pMerged' } (listOf \text{ parser}_{dig} \diamond \diamond \\ listOf \text{ parser}_{low} \diamond \diamond \\ listOf \text{ parser}_{upp})$$

The parser usage remains the same, but the result is a customized data structure. The *invert* function takes as many arguments as we have composed parsers.

5 Inverting parsers of different types

As we realized in the previous section, the type of the invert function for three combined parsers is:

$$invert :: l \rightarrow l \rightarrow l \rightarrow [l]$$

The three inputs have to be of the same type since they will be inserted in a list. However, parsers of different types can be combined, for instance the parser *dig_alp*, defined in the section 3.1. Then the *invert* function has to deal with inputs of different types. There are two ways to deal with this drawback. The most obvious one, is to return the inverted values in a data structure where the elements can be of different types, instead of a list. To achieve this solution, we need to replace the list insert function (*:*), used in the function *acceptNextVal*. We also have to specify the empty structure that *pFoldr* should return when the parser fails or reaches the end of input.

If we decide to return the inverted elements in a list, the user will have to supply functions, that homogenize all the parser results, in other words, all the inputs have to be transformed in values of the same type. This functions will be called *transformation* functions. The new *invert* function takes *n* *transformation* functions and *n* parser results. Each function homogenizes the respective parser result.

$$invert \ f_1 \dots f_n \ res_1 \dots res_n$$

Whenever we recover another element we have to homogenize it using the *transformation* function. All the *transformation* functions are stored in a nested pair. The function we should use depends on the parser that we used to recognize that element. We will encapsulate the type of each inverted value using a new type definition. In this data type, we will store two functions. The first function, *funcs* \rightarrow (*inv* \rightarrow *res*), knows how to choose the right *transformation function*, from a nested pair of functions. The second function, *ds* \rightarrow (*inv*, *ds*), is the function that takes the next inverted value from the parsing result.

$$\mathbf{data} \ Pair \ funcs \ res \ val = \\ \exists \ inv. \ Pair \ (funcs \rightarrow (inv \rightarrow res)) \ (val \rightarrow (inv, val))$$

With the second function *val* \rightarrow (*inv*, *val*) we get an inverted value of type *inv*. Then with the first function *funcs* \rightarrow (*inv* \rightarrow *res*), we get a function

that transform the inverted value in a value of type *res*. The \exists states that the inverted value, *inv*, can be of any type if we also provide a function that can transform that value to a value of type *res*. Using this data type we make the type of the inverted value transparent.

The new *invert* function takes *n transformation* functions and *n parser* results. Now we have to merge the *transformation* functions in a nested pair, as we do to the input of the *invert* function. Thus we need two *pack* functions.

As a practical example, we will redefine the *listOf* and *treeOf* functions. Now, instead of storing an *invert* function, we have to store a *Pair* with the function that is able to choose the *transformation* function for the respective inverted value, and the *invert* function itself.

listOf parser = ([], (($\lambda x l \rightarrow (x : l, \text{Pair } id \text{ getVal})$) \diamond *parser*), *id*, *id*, *id*)

treeOf parser = (*Leaf*, *f* \diamond *parser*, *id*, *id*, *id*)
where *f* *x* *t* = **let** (*res*, *getVal*) = *insert* *x* *t*
in (*res*, *Pair* *id* *getVal*)

The new definition of the combinator \diamond is extended with an extra *pack* function. This *pack* function will be used to pack the *transformation* functions, before being used by the *invert* function.

As you can realize, the two *pack* functions are similar. It will be convenient to have a single polymorphic function. That is possible, but it is out of scope of this paper:

(*e_a*, *p_a*, *punpack*, *ppack*, *ppack'*) \diamond (*e_b*, *p_b*, *qunpack*, *qpack*, *qpack'*) =
 ((*e_a*, *e_b*)
 , *comb_a* \diamond *p_a* \diamond *comb_b* \diamond *p_b*
 , *qunpack.uncurry.punpack*
 , *qpack.curry.ppack*
 , *qpack'.curry.ppack'*
)

The combined functions are adapted to the new type *Pair*. In *comb_a* first we apply *insert* to the left side of the pair. Then we update the *choose* function by applying the *ch* function to the left side of the pair of the *transformation* functions. The *invert* function takes a value from the left side, exactly how we did it before. The *comb_b*, has the same behavior but applies both functions to

the right side of the pair.

$$\begin{aligned}
\text{comb}_a \text{ insert } (values_a, values_b) &= ((newVal_a, values_b), \text{Pair choose invert}_a) \\
\text{where} & \\
(newVal_a, \text{Pair ch getVal}) &= \text{insert values}_a \\
\text{choose } (func_a, func_b) &= \text{ch func}_a \\
\text{invert}_a (inp_a, inp_b) &= \text{let } (e, rest) = \text{getVal } inp_a \\
&\quad \text{in } (e, (rest, inp_b)) \\
\text{comb}_b \text{ insert } (values_a, values_b) &= ((values_a, newVal_b), \text{Pair choose invert}_b) \\
\text{where} & \\
(newVal_b, \text{Pair ch getVal}) &= \text{insert values}_b \\
\text{choose } (func_a, func_b) &= \text{ch func}_b \\
\text{invert}_b (inp_a, inp_b) &= \text{let } (e, rest) = \text{getVal } inp_b \\
&\quad \text{in } (e, (inp_a, rest))
\end{aligned}$$

The *acceptNextVal* function uses the function *choose* to take the correct *transformation* function from the nested pair *funcs*, where all the *transformation* functions are stored. Then the chosen function is applied to the inverted value *e* and the result is inserted in the list of inverted values.

$$\begin{aligned}
\text{acceptNextVal } f (recognVals, getRest) &= (newValues, invert) \\
\text{where} & \\
(newValues, \text{Pair choose getVal}) &= f recognVals \\
\text{invert funcs inp} &= \text{let } (e, rest) = \text{getVal } inp \\
&\quad \text{in } (\text{choose funcs } e) : (\text{getRest funcs } rest)
\end{aligned}$$

The *acceptNextVal* is used by *pMerged* to compose the parsing results. By applying the two *pack* functions to *invert* we obtain an invert function that takes several inputs instead of two nested pairs.

$$\begin{aligned}
\text{sem 'pMerged' } (units, alternatives, \text{unpack}, \text{pack}, \text{pack}') &= \\
(\lambda(r, \text{invert}) \rightarrow (\text{unpack sem } r, \text{pack}' (\lambda val \rightarrow \text{pack } (\text{invert } val)))) & \\
\text{◇} & \\
\text{pFoldr acceptNextVal } (units, \text{const } (\text{const } [])) \text{ alternatives} &
\end{aligned}$$

Now we can combine parsers of different types. We can define a parser, that includes, for instance, the parser *parser_{int}* that recognizes and returns digits, as well as, the parser *parser_{low}* returns strings and the parser *parser_{upp}* that returns a balanced tree of characters.

$$\begin{aligned}
\text{int_low_upp} = (, ,) \text{ 'pMerged' } (\text{listOf parser}_{int} \text{ ◇} \\
\text{listOf parser}_{low} \text{ ◇} \\
\text{treeOf parser}_{upp})
\end{aligned}$$

The usage of the parser remains the same, but the result is a list of digits, a string and a tree of characters. To recover the order of the initial input, we supply to the invert function the functions that homogenize the input, for instance, the function *show*.

6 Example: Parsing Java syntax

Modern languages are liberal about the order of declarations. Taking as example a small subset of the java language specification, each **class** can be composed by fields, methods and constructors. These components can appear in any order.

We start by defining the abstract syntax for Java classes:

```
data Class = Class String Elements
data Elements = Elements [Field] [Constructor] [Method]
```

Assume we have three parsers $pField$, $pConstructor$ and $pMethod$, that parse fields, constructors and methods, respectively.

Then we use the $pMerged$ function to build the final parser $pClass$. Each Java field, constructor and method is stored in a list, using the function $listOf$.

```
pClass = Class <| pToken "class" <| identifier <|
braces (Elements 'pMerged' (
    listOf pField <|
    listOf pConstructor <|
    listOf pMethod
))

braces p = (λ_ v _ → v) <| pToken "{" <| p <| pToken "}"
```

To recover the initial Java code we only need to define a pretty-printing for class and use the resulting $invert$ function. This function will return the input in its original order.

```
ppClass invert (Class name (Elements fs cs ms)) =
    "class " ++ name ++
    ppBraces (invert ppField ppConstructor ppMethod fs cs ms)
```

7 Conclusion

We have shown how to define a *NOP* parser by combining several parsers. We extended the parser combinator with the inverse functionality.

We tried to make the interface of the combinators and the invert function as simple as possible. A user can easily make use of this combinator to obtain a *NOP* parser and the invert function, with the advantage of having the parsing result in a customized data type.

We have shown how to use existentially quantified data types to customize the type of the parser result. The invert function can be used to generate properly located error messages or even to repair errors found at parsing time.

The functions presented in this project have been incorporated into the "UU_Parsing" library. This library is used by a large community of functional programmers worldwide. Thus the results of this project will also have an impact in that community and we do hope to contribute to make programming and writing parsers a more effective, easy and grateful task.

References

- [1] A.I. Baars, A. Löh, and S.D. Swierstra. Parsing permutation phrases. In R. Hinze, editor, *Proceedings of the 2001 ACM SIGPLAN Haskell Workshop*, pages 171–182. Elsevier, 2001.
- [2] J. Fokker. Functional parsers. In J.T. Jeuring and H.J.M. Meijer, editors, *Advanced Functional Programming*, number 925 in LNCS, pages 1–52, 1995.
- [3] Simon Peyton-Jones, John (eds.) Hughes, et al. Report on the programming language Haskell 98. February 1998.
- [4] Doaitse Swierstra and Luc Duponcheel. Deterministic, error correcting combinator parsers. In *Advanced Functional Programming, Second International Spring School*, volume 1129 of LNCS, pages 184–207. Springer-Verlag, 1996.
- [5] Philip Wadler. How to replace failure with a list of successes. In *Functional Programming Languages and Computer Architecture*, volume 201 of LNCS, pages 113–128. Springer-Verlag, 1985.

A Parsing a permutation phrase with inverse functionality

In [1] is shown how to represent a permutation phrase in a tree by expanding and factorizing its elements. Each path from the root to a leaf in the tree represents a particular permutation. Permutations with a common prefix share the same subtree, hence the number of choices in each node is limited by the number of permutable elements.

The data type *Perms* represents our tree. If all elements stored in a tree are optional then their default values are stored in *defaults*, otherwise *defaults* is *Nothing*. The parser stored in each Branch is not allowed to derive the empty string.

```
data Perms a = Choice{ defaults :: (Maybe a), branches :: [Branch a] }
data Branch a =  $\forall x$ . Br (Parser (x, GetVal)) (Perms (x  $\rightarrow$  a))
```

The first parser added to the tree, will have an invert function that takes the first element of the input list. The *n* parser, where *n* denotes the number of parsers added to the tree plus one, will take the element in the position number *n* from the input list.

```
data GetVal = GetVal ( $\forall a$ . [a]  $\rightarrow$  a)
```

The data type *GetVal* keeps record of the function that takes the *n* element of a list.

```
data Undo a = Undo GetVal (Perms a)
```

By keeping this information outside the tree, we avoid to replicate it through all the branches.

All the parsers of the permutation phrase has to be added to the tree. Each path of the tree represents a parser of a possible permutation. The function *permute* follows the correct path to obtain the pretended parser. The new parser returns the parsing result and the *invert* function.

```
permute :: Undo a  $\rightarrow$  Parser (a, [b]  $\rightarrow$  ([b], [b]))
permute (Undo g (Choice def bs)) = foldr (<math>\langle \diamond \rangle</math>) exit (map pars bs)
where exit = case def of
  Nothing  $\rightarrow$  fail
  Just x  $\rightarrow$  succeed (x,  $\lambda inp \rightarrow$  (inp, []))
  pars (Br p perm) = ( $\lambda(x, GetVal u) (f, l) \rightarrow$  (f x, (nextVal u).l))  $\langle \diamond \rangle$ 
    p  $\langle \diamond \rangle$  permute (Undo g perm)
```

```
nextVal getVal (vals, inverted) = (vals, (getVal vals) : inverted)
```

The empty tree stores the initial *invert* function, (*GetVal head*).

```
empty :: a  $\rightarrow$  Undo a
empty x = Undo (GetVal head) (Choice (Just x) [])
```

