

Parsing Permutation Phrases

Arthur Baars, Andres Löh, Doaitse Swierstra

Institute of Information and Computing Science
Utrecht University

e-mail: {arthurb, andres, doaitse}@cs.uu.nl

September 2, 2001

Permutation phrases

All the following describe the same tag:

```
  
  

```

- Attributes of XML tags
- Specifiers in C-declarations
 - `static int signed foo;`
 - `int static signed foo;`
- Derived read-function for Haskell's record syntax
- Citation fields in BIB_{TEX} entries
-

Previous work

Extending context-free grammars with permutation phrases

Robert D. Cameron:

- Extension of EBNF notation for permutation phrases
 - $\ll A \parallel B \parallel C \gg$
- Pseudo-code algorithm
 - LL(1) parsing
 - Parser generators
 - Ignores reordering and type of different constituents

Interface of our parser combinators

We will present our permutation combinators in the context of the latest version of the UU_Parsing library:

infixl 5 \triangleleft

infixl 6 \triangleleft^*

infix 7 $\triangleleft^\$$

class *IsParser* *p* **where**

pFail :: *p a*

pSucceed :: *a* \rightarrow *p a*

pSym :: *Char* \rightarrow *p Char*

(\triangleleft^*) :: *p (a* \rightarrow *b)* \rightarrow *p a* \rightarrow *p b*

(\triangleleft) :: *p a* \rightarrow *p a* \rightarrow *p a*

$(\triangleleft^\$)$:: *(a* \rightarrow *b)* \rightarrow *p a* \rightarrow *p b*

f $\triangleleft^\$$ *p* = *pSucceed f* \triangleleft^* *p*

parse :: *p a* \rightarrow *String* \rightarrow *Result a*

How to use parser combinators

```
src = "/icons/li-dsol.gif"  
width = 190
```

$$\begin{aligned} pField &:: \text{IsParser } p \Rightarrow \text{String} \rightarrow p \ a \rightarrow p \ a \\ pField \ fld \ p &= (\lambda_ _ \text{val} \rightarrow \text{val}) \ \langle \$ \rangle \ pTok \ fld \ \langle * \rangle \ pSym \ '=' \ \langle * \rangle \ p \end{aligned}$$

Adding parentheses according to associativity and priority of $\langle \$ \rangle$ and $\langle * \rangle$:

$$pField \ fld \ p = (((((\lambda_ _ \text{val} \rightarrow \text{val}) \ \langle \$ \rangle \ pTok \ fld) \ \langle * \rangle \ pSym \ '=') \ \langle * \rangle \ p)$$
$$\begin{aligned} pMaybe &:: \text{IsParser } p \Rightarrow p \ a \rightarrow p \ (\text{Maybe } a) \\ pMaybe \ p &= \text{Just } \langle \$ \rangle \ p \\ &\langle \rangle \ pSucceed \ \text{Nothing} \end{aligned}$$

```
UU_Parsing_Demo> parse (pMaybe (pField "width" pLength)) "width_=_190"
```

```
Result: Just 190
```

```
UU_Parsing_Demo> parse (pMaybe (pField "width" pLength)) ""
```

```
Result: Nothing
```

How to use the new combinators

```

```

```
data XHTML = Img { src      :: URI
                  , alt      :: Text
                  , width    :: Maybe Length
                  , height   :: Maybe Length
                  }
```

```
pImgTag      :: Parser XHTML
```

```
pImgTag      = ( $\lambda\_ \_ x \_ \rightarrow x$ ) <$> pTok "<" <*> pTok "img" <*> attrs <*> pTok ">"
```

```
  where attrs = pPerms (Img <$> pField "src" pURI
                       <*> pField "alt" pText
                       <*> pMaybe (pField "width" pLength)
                       <*> pMaybe (pField "height" pLength)
                       )
```

Interface of new parser combinators

infix 6 $\langle * \rangle$, $\langle \langle * \rangle \rangle$

infix 7 $\langle \$ \rangle$, $\langle \langle \$ \rangle \rangle$

$\langle * \rangle \quad :: \text{IsParser } p \Rightarrow \quad p (a \rightarrow b) \rightarrow p a \rightarrow \quad p b$

$\langle \langle * \rangle \rangle \quad :: \text{IsParser } p \Rightarrow \text{Perms } p (a \rightarrow b) \rightarrow p a \rightarrow \text{Perms } p b$

$\langle \$ \rangle \quad :: \text{IsParser } p \Rightarrow (a \rightarrow b) \rightarrow p a \rightarrow \quad p b$

$\langle \langle \$ \rangle \rangle \quad :: \text{IsParser } p \Rightarrow (a \rightarrow b) \rightarrow p a \rightarrow \text{Perms } p b$

$p\text{Perms} \quad :: \text{IsParser } p \Rightarrow \text{Perms } p a \rightarrow p a$

Toy Example

$S ::= \ll a^* \parallel b \parallel c? \gg$

$pS :: \text{Parser } (\text{String}, \text{Char}, \text{Char})$
 $pS = pPerms ((,,) \langle\langle \$ \rangle\rangle pList (pSym 'a')$
 $\langle\langle * \rangle\rangle pSym 'b'$
 $\langle\langle ? \rangle\rangle pOptSym 'c'$
 $)$

$pOptSym :: \text{Char} \rightarrow \text{Parser Char}$
 $pOptSym x = pSym x \langle\langle \rangle\rangle pSucceed '_'$

UU_Parsing_Demo> *parse* pS "caab"

Result:

("aa", 'b', 'c')

Toy Example (continued)

$S ::= \langle\langle a^* \parallel b \parallel c? \rangle\rangle$

```
UU_Parsing_Demo> parse pS ""
```

Symbol 'b' inserted at end of file; 'b' or 'c' or ('a')* expected.

Result:

```
("", 'b', '_')
```

```
UU_Parsing_Demo> parse pS "abd"
```

Errors:

Symbol 'd' at end of file deleted; 'c' or eof expected.

Result:

```
("a", 'b', '_')
```

Requirements

- Easy to construct parsers for permutation phrases
- Efficient
- Inheritance of error-reporting/repairing facilities
- Constituents of the permutation phrase may have different types!!!
- Constituents of the permutation phrase may be optional (empty)

Implementation Idea

The permutation phrase:

$$S ::= \ll a \parallel b \parallel c \gg$$

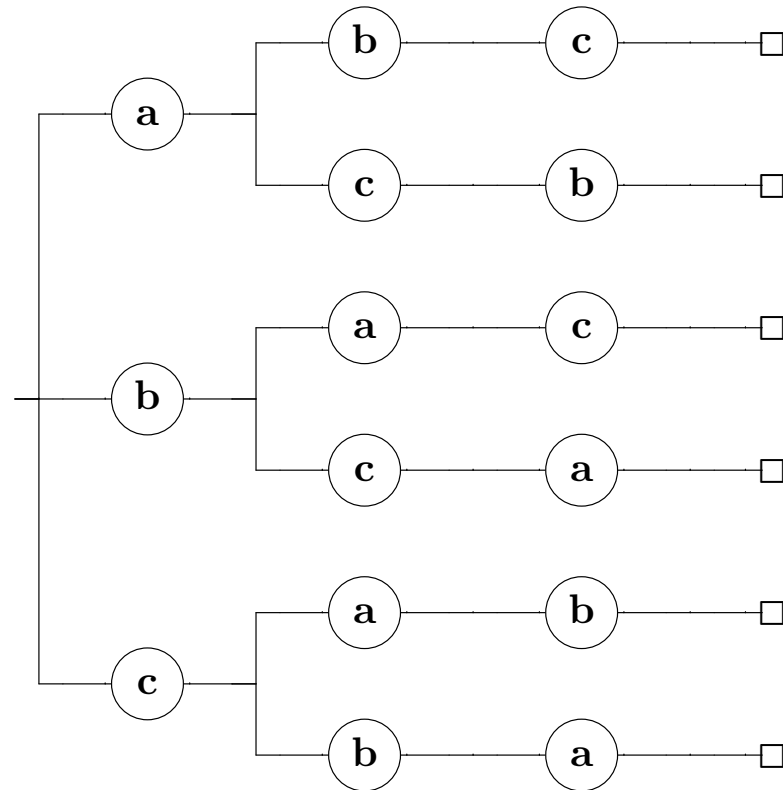
is formally equivalent to:

$$\begin{array}{l} S ::= abc \\ | acb \\ | bac \\ | bca \\ | cab \\ | cba \end{array}$$

Implementation Idea (continued)

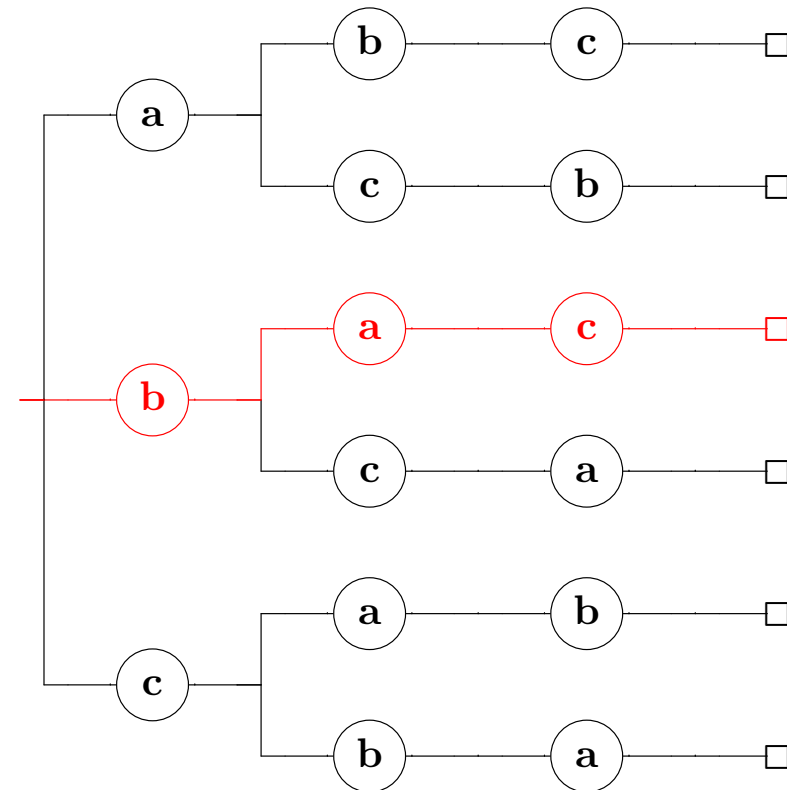
$$\begin{aligned} S & ::= a(bc|cb) \\ & | b(ac|ca) \\ & | c(ab|ba) \end{aligned}$$

Consider parsing the string “bac”



Implementation Idea (continued)

- When recognizing a permutation a parser follows a path through the permutation tree from the root to one of the exits.
- Permutation tree has a potential size of order $n!$, where n is the number of permutable constituents.
- Lazy evaluation ensures that only those parts of the tree are evaluated that correspond to the path followed by the parser.



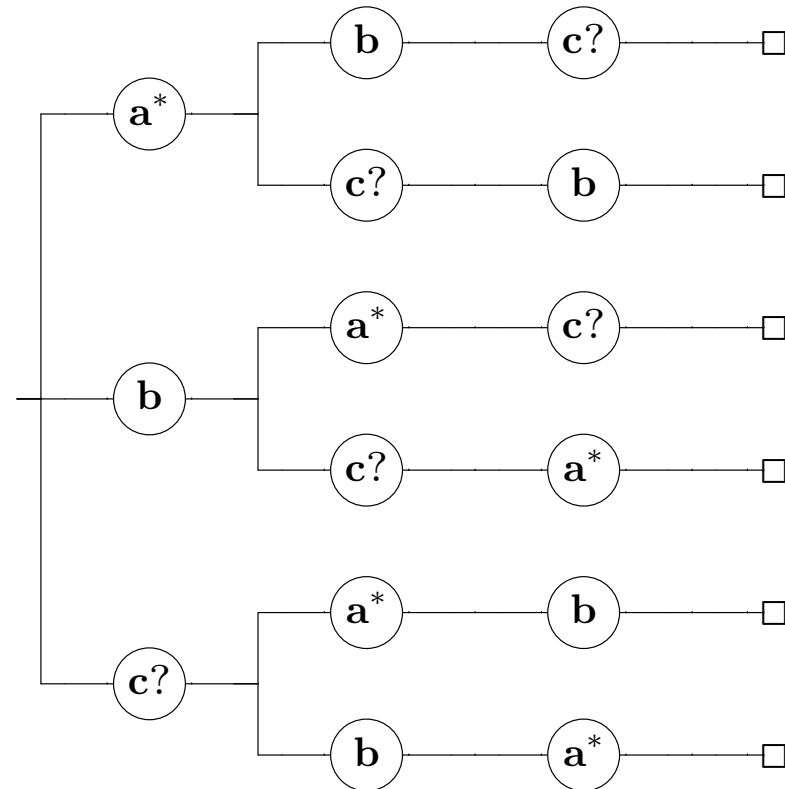
Complications caused by possibly empty elements

$S ::= \ll a^* \parallel b \parallel c? \gg$

$S ::= a^*(bc?|c?b)$
 $\quad | b(a^*c?|c?a^*)$
 $\quad | c?(a^*b|ba^*)$

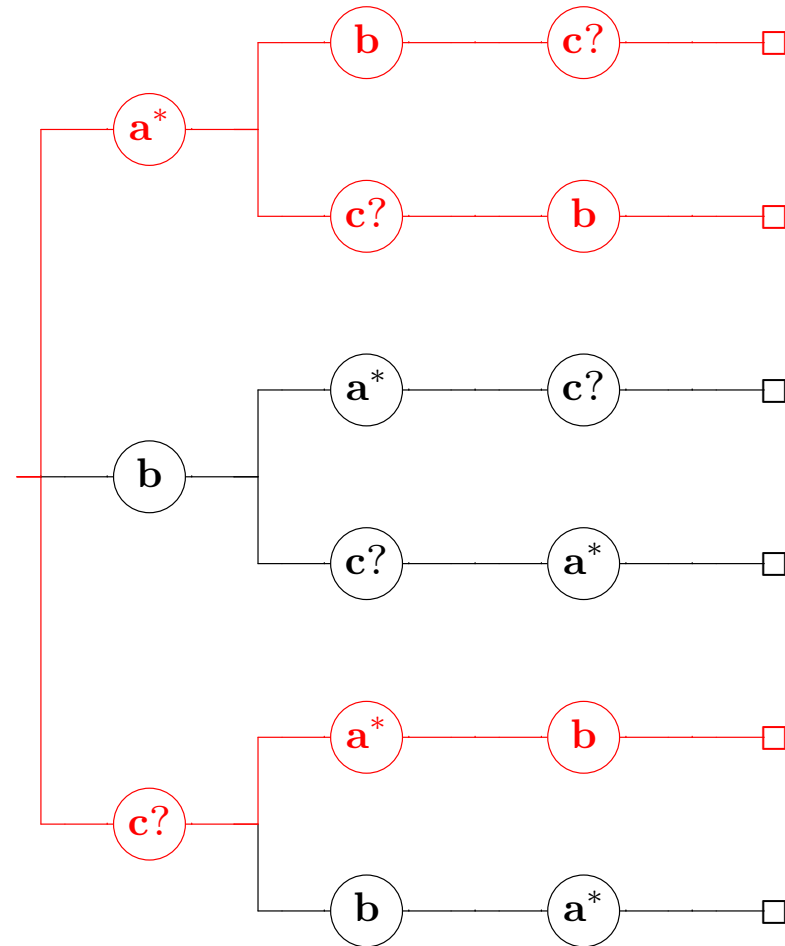
Elements with a ? or a * can be empty.

Problem: “aab” is ambiguous.



Complications caused by possibly empty elements (continued)

There are three different paths for “aab”, because the empty “c?” can be derived at different positions:



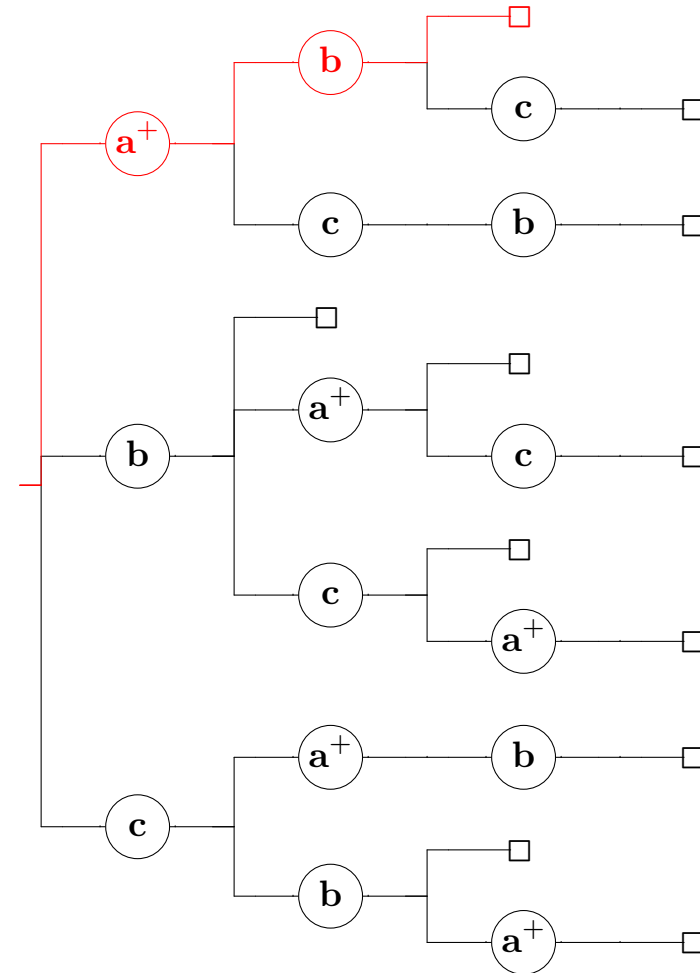
Disambiguating: splitting constituents in empty and non-empty part

Solution: optional elements that are not present can only be derived at the end of a parse.

This means: no nodes can derive the empty string, and additional exits (leaves) must be added.

splitParser takes care of splitting an optional element (parser that can derive the empty string) in its empty part and its non-empty part.

$splitParser :: p\ a \rightarrow (Maybe\ a, Maybe\ (p\ a))$

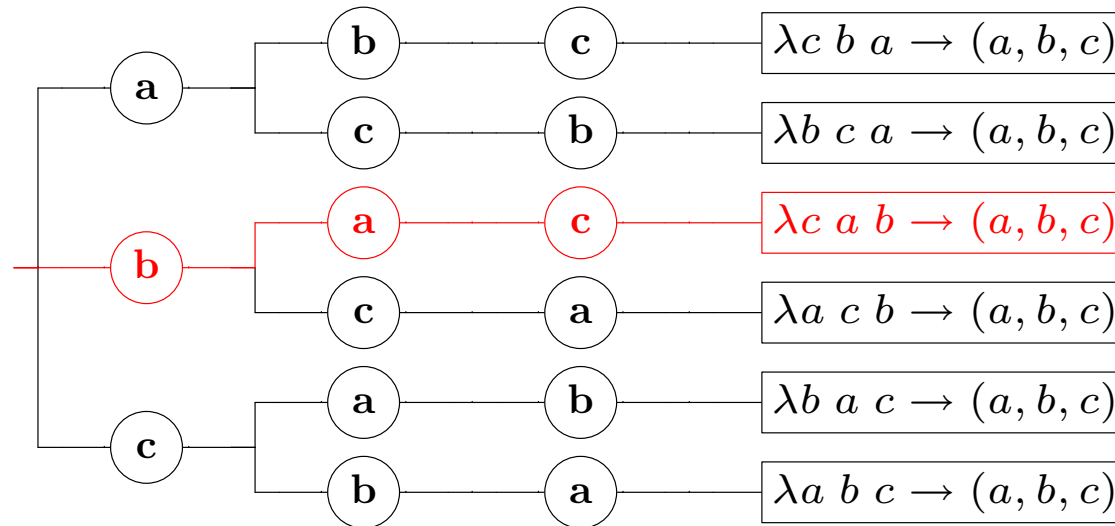


Semantic functions are stored at the exits

$pABC \quad :: \text{Perms } (\text{Char}, \text{Char}, \text{Char})$

$pABC \quad = p\text{Perms } ((,,) \langle\langle \$ \rangle\rangle p\text{Sym } 'a' \langle\langle * \rangle\rangle p\text{Sym } 'b' \langle\langle * \rangle\rangle p\text{Sym } 'c')$

- In every exit a version of the semantic function is stored that that arranges the arguments in a uniform order.

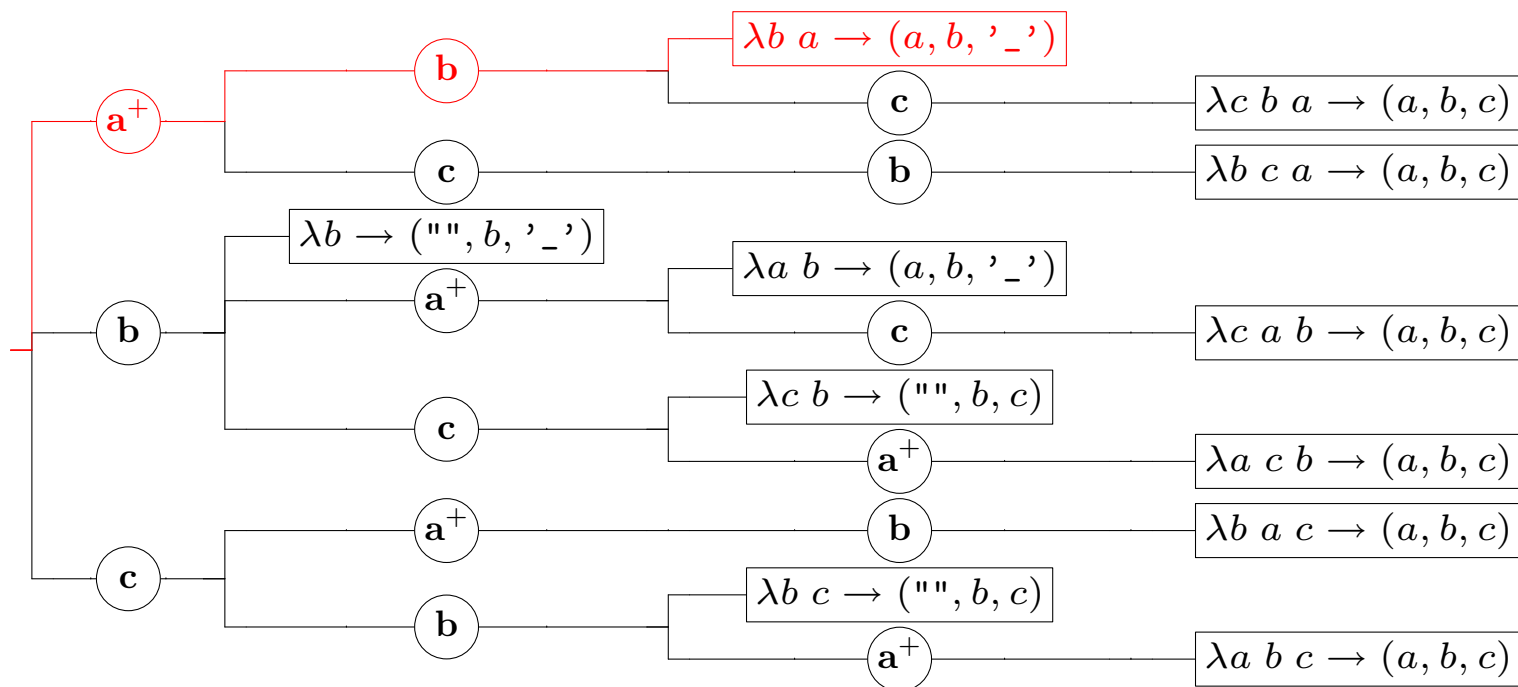


Toy example revisited

$pS \quad :: \text{Parser } (\text{String}, \text{Char}, \text{Char})$

$pS \quad = \text{pPerms } ((,,) \ll \$ \gg \text{pList } (\text{pSym } 'a') \ll * \gg \text{pSym } 'b' \ll * \gg \text{pOptSym } 'c')$

$\text{pOptSym } x = \text{pSym } x \triangleleft \text{pSucceed } '_'$



Implementation: *Perms* datatype

data *Perms* *p* *a* = *Choice* { *exit* :: *Maybe* *a*, *branches* :: [*Branch* *p* *a*] }

data *Branch* *p* *a* = $\forall x. Br (p\ x) (Perms\ p\ (x \rightarrow a))$

Note the use of an existential type to hide the specific order in which the elements occur.

Definitions of *fmap* for *Perms* *p* and *Branch* *p*:

instance *Functor* (*Perms* *p*) **where**

$$fmap\ v_{a \rightarrow b}\ (Choice\ e_a\ bs_a) = Choice\ (fmap\ v_{a \rightarrow b}\ e_a) \\ [fmap\ v_{a \rightarrow b}\ br_a \mid br_a \leftarrow bs_a]$$

instance *Functor* (*Branch* *p*) **where**

$$fmap\ v_{a \rightarrow b}\ (Br\ p_x\ t_{x \rightarrow a}) = Br\ p_x\ (fmap\ (v_{a \rightarrow b} \circ)\ t_{x \rightarrow a})$$

Implementation: transforming *Perms* into a parser

$pPerms\ t_a$ $=\ foldr\ (\triangleleft)\ empty\ nonempties$
 where $empty$ $=\ \mathbf{case\ } exit\ t_a\ \mathbf{of}$
 $Nothing \rightarrow pFail$
 $Just\ v_a \rightarrow pSucceed\ v_a$
 $nonempties$ $=\ [flip\ (\$)\ \triangleleft\$ p_x\ \triangleleft* pPerms\ t_{x \rightarrow a}$
 $| Br\ p_x\ t_{x \rightarrow a} \leftarrow branches\ t_a$
 $]$

Building the *Perms* structure

$succeedPerms \quad :: \quad a \rightarrow Perms \ p \ a$
 $succeedPerms \ x \quad = \quad Choice \ (Just \ x) \ []$

$(\langle\langle \$ \rangle\rangle)$ $:: \quad IsParser \ p \Rightarrow (a \rightarrow b) \rightarrow p \ a \rightarrow Perms \ p \ b$
 $f \langle\langle \$ \rangle\rangle p \quad = \quad succeedPerms \ f \ \langle\langle * \rangle\rangle \ p$

$(\langle\langle * \rangle\rangle)$ $:: \quad IsParser \ p \Rightarrow Perms \ p \ (a \rightarrow b) \rightarrow p \ a \rightarrow Perms \ p \ b$
 $perms \ \langle\langle * \rangle\rangle \ p \quad = \quad add \ (splitParser \ p) \ perms$

Example of usage:

$pS \quad :: \quad Parser \ (String, \ Char, \ Char)$
 $pS \quad = \quad pPerms \ ((, ,) \ \langle\langle \$ \rangle\rangle \ pList \ (pSym \ 'a' \ \langle\langle * \rangle\rangle \ pSym \ 'b' \ \langle\langle * \rangle\rangle \ pOptSym \ 'c' \)$

Building the *Perms* structure (continued)

empty part	non-empty part	
<i>Nothing</i>	<i>Nothing</i>	<i>pFail</i>
<i>Just</i>	<i>Nothing</i>	<i>pSucceed</i>
<i>Nothing</i>	<i>Just</i>	required element
<i>Just</i>	<i>Just</i>	optional element

$add :: (Maybe a, Maybe (p a)) \rightarrow Perms p (a \rightarrow b) \rightarrow Perms p b$
 $add\ sp_a\ t_{a \rightarrow b}@(Choice\ e_{a \rightarrow b}\ bs_{a \rightarrow b}) = \mathbf{case\ } sp_a\ \mathbf{of}$
 $(Nothing, Nothing) \rightarrow Choice\ Nothing\ []$
 $(Just\ v_a, Nothing) \rightarrow fmap\ (\$v_a)\ t_{a \rightarrow b}$
 $(Nothing, Just\ p_a) \rightarrow Choice\ Nothing\ (insert\ p_a)$
 $(Just\ v_a, Just\ p_a) \rightarrow Choice\ (fmap\ (\$v_a)\ e_{a \rightarrow b})\ (insert\ p_a)$
where $insert\ p_a = Br\ p_a\ t_{a \rightarrow b} : [Br\ p_x\ (add\ sp_a\ t_{a \rightarrow x \rightarrow b})$
 $\quad | Br\ p_x\ t_{x \rightarrow a \rightarrow b} \leftarrow bs_{a \rightarrow b}$
 $\quad , \mathbf{let}\ t_{a \rightarrow x \rightarrow b} = fmap\ flip\ t_{x \rightarrow a \rightarrow b}$
 $\quad]$

Conclusions

- Existential types
 - Reordering of constituents is type-safe and hidden from user
- Lazy evaluation
 - Only small part of intermediate data structure is evaluated
- No need for a second phase: errors are reported by underlying parsing library
- Can be made to work with (almost) any parser combinator library
- http://www.cs.uu.nl/groups/ST/Software/UU_Parsing/