

# Derivation of a Typed Functional LR Parser

Ralf Hinze ([ralf@informatik.uni-bonn.de](mailto:ralf@informatik.uni-bonn.de))

*Institut für Informatik III, Universität Bonn, Römerstraße 164, 53117 Bonn, Germany.*

Ross Paterson ([ross@soi.city.ac.uk](mailto:ross@soi.city.ac.uk))

*Department of Computing, City University, London EC1V 0HB, United Kingdom.*

**Abstract.** This paper describes a purely functional implementation of LR parsing. We formally derive our parsers in a series of steps starting from the inverse of printing. In contrast to traditional implementations of LR parsing, the resulting parsers are fully typed, stackless and table-free. The parsing functions pursue alternatives in parallel with each alternative represented by a continuation argument. The direct implementation presents many opportunities for optimization and initial measurements show excellent performance in comparison with conventional table-driven parsers.

**Keywords:** continuations, derivation, LL- and LR parsing, efficiency.

## 1. Introduction

Variants of LR parsing (Knu65) have long been the favoured technique for mechanically generated parsers. LR grammars are expressive, allowing a convenient description of most common programming language constructs, and many tools exist to generate efficient parsers from them.

Typical implementations reflect the origin of these parsers as deterministic pushdown automata: there is a table of actions for each state and input symbol, and a stack holding saved values and states. In this paper, we shall derive an alternative implementation: the parsers developed here consist of mutually recursive functions, each taking a number of continuation functions. To give an idea of our destination, Figure 1 lists a complete parser for a tiny expression language. The states of conventional parsers have become functions with continuation arguments, and the stack has disappeared: values discovered during parsing are immediately passed to the continuations. Another salient feature of our parsers is that they are fully typed, without the need for a union type for stack elements as in many other implementations.

We shall explain our parsers by deriving them in a series of program transformation steps. After introducing the notations we shall be using in Section 2, the development is as follows:

- In Section 3, we begin with a direct specification of parsing, as the inverse of flattening a parse tree. Parsing is therefore a non-deterministic, or set-valued, function. It is convenient to work with

```

data Expr      = Ident String | Apply Expr Expr
reduceid      :: (Expr → Parser r) → (String → Parser r)
reduceid g n   = g (Ident n)
reduceE(E)    :: (Expr → Parser r) → (Expr → Expr → Parser r)
reduceE(E) g f x = g (Apply f x)
data Token     = IDENT String | LPAREN | RPAREN | EOF
type Parser r  = [Token] → Maybe r
parseE        :: Parser Expr
parseE        = state1 (λe tr → return e)
state1        :: (Expr → Parser r) → (Parser r)
state1 k ts    = case ts of IDENT n : tr → reduceid g n tr
                    -                → fail "syntax error"
    where g e   = state2 (reduceE(E) g e) (k e)
state2        :: (Expr → Parser r) → (Parser r) → (Parser r)
state2 k1 k2 ts = case ts of LPAREN : tr → state3 k1 tr
                    EOF : tr    → k2 tr
                    -                → fail "syntax error"
state3        :: (Expr → Parser r) → (Parser r)
state3 k ts    = case ts of IDENT n : tr → reduceid g n tr
                    -                → fail "syntax error"
    where g e   = state4 (reduceE(E) g e) (k e)
state4        :: (Expr → Parser r) → (Parser r) → (Parser r)
state4 k1 k2 ts = case ts of LPAREN : tr → state3 k1 tr
                    RPAREN : tr → k2 tr
                    -                → fail "syntax error"

```

Figure 1. Parser for the grammar  $E \rightarrow \text{id} \mid E ( E )$ .

such functions through most of our development, addressing the deterministic special case towards the end.

- The next step, in Section 4, is to derive a simple recursive characterization of the parsing functions. This form corresponds to continuation-style recursive descent parsers, and is thus directly executable only if the grammar is not left-recursive.
- Left recursion is removed in Section 5 by analysing these functions as sets of sequences of primitive steps, given by a left-linear grammar, and using the standard technique to obtain an equivalent right-linear grammar. The resulting program corresponds to a non-deterministic LR parser.
- Finally, in Section 6, we remove non-determinism by collecting sets of parsing functions, so that they parse alternatives in paral-

lel, obtaining deterministic LR parsers. We also consider various optimizations. Some of these have counterparts for conventional parsers; some are specific to the functional representation.

We discuss previous derivations of bottom-up functional parsers in Section 7, and compare them with the parsers derived here.

Experiments with a range of grammars, reported in Section 8, show that our parsers compare well to conventional table-based ones.

Finally, Section 9 concludes.

## 2. Notation

This section briefly introduces the notation used in the subsequent sections. The reader may wish to skim through the material and come back later for reference.

*Sets* As noted in the introduction we shall derive an implementation of LR parsing in a series of steps, each of which yields an executable program. For concreteness, the derivation and the resulting programs are presented in a Haskell-like notation; for information on Haskell, see (Pey03). Although we use Haskell as a vehicle, we are interpreting it in the category of sets and total functions so that arbitrary recursion is unavailable. At the end of the development, when all non-determinism is removed, we shall switch to the conventional interpretation of Haskell in the category of complete partial orders and continuous functions.

Furthermore, we shall make heavy use of set comprehension notation. As an example, `prefixes` is a function that yields the set of all prefixes of a given list.

$$\begin{aligned} \text{prefixes} &:: [\mathbf{a}] \rightarrow \mathbb{P} [\mathbf{a}] \\ \text{prefixes } x &= \{ p \mid p \# s = x \} \end{aligned}$$

Here  $[\mathbf{a}]$  is the type of lists of elements of type  $\mathbf{a}$ ; likewise,  $\mathbb{P} \mathbf{a}$  is the type of sets of elements of type  $\mathbf{a}$ ; ‘ $\#$ ’ denotes list concatenation. Note that the variables  $p$  and  $s$  in the body of the comprehension are implicitly existentially quantified so that the definition is shorthand for the more verbose

$$\text{prefixes } x = \{ p \mid \exists p :: [\mathbf{a}]. \exists s :: [\mathbf{a}]. p \# s = x \}$$

For the derivation we prefer the shorter form, so to reduce clutter we adopt the convention of (Hin02) that the variables to the left of ‘ $=$ ’ and ‘ $\in$ ’ are always implicitly existentially quantified.

*Context-free grammars* We shall use standard notation for context-free grammars with one minor exception: we allow a set of start symbols. Formally, a *context-free grammar*  $(\mathbf{T}, \mathbf{N}, \mathbf{S}, \mathbf{P})$  consists of

- disjoint sets  $\mathbf{T}$  and  $\mathbf{N}$ , called respectively *terminal* and *nonterminal symbols*, together comprising the set  $\mathbf{V} = \mathbf{T} \cup \mathbf{N}$  of *symbols*,
- a set  $\mathbf{S} \subseteq \mathbf{N}$  of *start symbols*, and
- a set  $\mathbf{P} \subseteq \mathbf{N} \times \mathbf{V}^*$  of *productions*.

We shall use the following conventions for elements of these sets:

$$\begin{array}{ll} a, b, c, d & \in \mathbf{T} & w & \in \mathbf{T}^* \\ A, B, C, D & \in \mathbf{N} & & \\ X & \in \mathbf{V} & \alpha, \beta, \gamma & \in \mathbf{V}^* \end{array}$$

Elements of  $\mathbf{P}$  are written  $A \rightarrow \alpha$ .

*Example 1.* Consider the simple expression grammar (ASU86, p222) with

$$\begin{array}{l} \mathbf{T} = \{ +, *, (, ), \text{id} \} \\ \mathbf{N} = \{ E, T, F \} \\ \mathbf{S} = \{ E \} \end{array}$$

and the following productions:

$$\begin{array}{l} E \rightarrow E + T \\ E \rightarrow T \\ T \rightarrow T * F \\ T \rightarrow F \\ F \rightarrow ( E ) \\ F \rightarrow \text{id} \end{array}$$

We shall use this grammar as a running example. □

For any grammar, we can define an *augmented grammar* (we tacitly assume that  $\$ \notin \mathbf{V}$ ):

$$\begin{array}{l} \mathbf{T}^\dagger = \mathbf{T} \cup \{ \$ \} \\ \mathbf{N}^\dagger = \mathbf{N} \cup \mathbf{S}^\dagger \\ \mathbf{S}^\dagger = \{ S^\dagger \mid S \in \mathbf{S} \} \\ \mathbf{P}^\dagger = \mathbf{P} \cup \{ S^\dagger \rightarrow S\$ \mid S \in \mathbf{S} \} \end{array}$$

In an augmented grammar the start symbols do not appear on the right-hand side of an production. Furthermore, the end of input is signalled by the special symbol '\$'.

*Linear grammars* A *left-linear grammar* is one in which all productions have the form  $A \rightarrow Bw$  or  $A \rightarrow w$ . Similarly, a *right-linear grammar* has only productions of the forms  $A \rightarrow wB$  and  $A \rightarrow w$ .

Any left-linear grammar  $(\mathbf{T}, \{S\} \uplus \mathbf{N}, \{S\}, \mathbf{P})$  with a single start symbol  $S$  that does not appear on the right-hand side of a production may be transformed into an equivalent right-linear grammar  $(\mathbf{T}, \{S\} \uplus \bar{\mathbf{N}}, \{S\}, \bar{\mathbf{P}})$ , or *vice versa*, by the rules:

$$\begin{array}{llll} \mathbf{P} \ni S \rightarrow Aw & \iff & \bar{A} \rightarrow w & \in \bar{\mathbf{P}} \\ S \rightarrow w & \iff & S \rightarrow w & \\ A \rightarrow Bw & \iff & \bar{B} \rightarrow w\bar{A} & \\ A \rightarrow w & \iff & S \rightarrow w\bar{A} & \end{array}$$

*Items* An *item* is a production plus a position in the right-hand side, written  $A \rightarrow \alpha \cdot \beta$  where the dot marks the position. *Kernel items* of an augmented grammar  $(\mathbf{T}^\dagger, \mathbf{N}^\dagger, \mathbf{S}^\dagger, \mathbf{P}^\dagger)$  are those of the forms  $S^\dagger \rightarrow \cdot S^\dagger$  or  $A \rightarrow \alpha \cdot \beta$  with  $\alpha \neq \epsilon$ . We let  $i$  and  $j$  denote arbitrary items, and use  $q$  for kernel items.

We also define a predict relation ‘ $\triangleright$ ’ on items by  $(A \rightarrow \alpha \cdot C\beta) \triangleright (C \rightarrow \cdot \gamma)$ , where  $C \rightarrow \gamma \in \mathbf{P}^\dagger$ , with transitive closure  $\triangleright^+$  and reflexive-transitive closure  $\triangleright^*$ . Furthermore, let  $q^+ = \{i \mid q \triangleright^+ i\}$  and  $q^* = \{i \mid q \triangleright^* i\}$ . Note that  $q^+ = q^* \setminus \{q\}$  if and only if  $q$  is a kernel item. For example, the equation does not hold for  $A \rightarrow \cdot Aa$ .

Finally, for any set  $Q$  of items, we define

$$\begin{aligned} Q/X &= \{A \rightarrow \alpha \cdot X\beta \in Q\} \\ Q+X &= \{A \rightarrow \alpha X \cdot \beta \mid A \rightarrow \alpha \cdot X\beta \in Q\} \\ \text{next}(Q) &= \{X \mid A \rightarrow \alpha \cdot X\beta \in Q\} \end{aligned}$$

### 3. A specification of parsing

#### 3.1. ATTRIBUTES

In most parsing applications, we wish to collect information about the string parsed. The standard approach is to define for each symbol  $X$  a sequence of data types  $\bar{X}$  of *attributes* computed for each parse of the symbol. (We consider only synthesized attributes here.) It would suffice to have a single attribute type for each symbol, representing multiple attributes by tuples and no attributes by a unit type. However, we choose to retain the sequences of types because they are more efficient, and do not significantly complicate the description as we will need to

deal with sequences of types anyway. We adopt the notational device of extending the function type to sequences of types, allowing

$$t_1 \cdots t_n \rightarrow t$$

as an alternative form of

$$t_1 \rightarrow \cdots \rightarrow t_n \rightarrow t$$

Once we have defined  $\bar{X}$  for each symbol  $X$ , we can form  $\bar{\alpha}$  by concatenating these sequences of types for each symbol in  $\alpha$ . We shall also write  $v_\alpha$  for a sequence of variables with types  $\bar{\alpha}$ .

### 3.2. TOKENS

We assume for each terminal symbol  $a$ , a string of attribute types  $\bar{a}$ , with  $\bar{\epsilon} = \epsilon$ . The combination of a terminal symbol with attribute values is conventionally called a *token* (AU72). We can represent these with a Haskell data type as follows:

$$\mathbf{data} \text{Token} = \sum_{c \in \mathbf{T}^\dagger} \mathbf{T}_c \bar{c}$$

For example, a `Token` type and a suitable lexical analyser for the expression grammar, Example 1, are given in Figure 2. In this case, some symbols have a single attribute, while others have none. Note that ‘+’ and ‘\*’ have been generalized to additive and multiplicative operator families whose attribute specifies the concrete operator at hand. The token EOF corresponds to  $\mathbf{T}_\S$ .

### 3.3. PARSE TREES

In general, productions would be annotated with rules for computing the attributes of the left side from those of the right side. Our aim here, though, is to derive parsing as an invertible process, so we shall use parse trees as attributes, thus ensuring that no information about the parse is lost. Hence we define for each nonterminal symbol  $A$ , a parse tree type  $\bar{A}$  as follows

$$\mathbf{data} \bar{A} = \sum_{(A \rightarrow \alpha) \in \mathbf{P}} \mathbf{Rule}_{A \rightarrow \alpha} \bar{\alpha}$$

Note that  $\mathbf{Rule}_{A \rightarrow \alpha}$  has type  $\bar{\alpha} \rightarrow \bar{A}$  demonstrating that the parse tree constructor is the inverse of a production. As an exception, in the case of the start symbols  $S^\dagger$  of an augmented grammar, we define  $\bar{S}^\dagger = \bar{S}$ . For example, parse tree types for the expression grammar appear in Figure 3.

Parse trees are an initial algebra of the signature determined by the grammar. The attributes that are usually associated with nonterminal

```

data Token = IDENT String
           | ADDOP AddOp           -- '+' and '-'
           | MULOP MulOp           -- '*' and '/'
           | LPAREN
           | RPAREN
           | EOF                   -- = T§
data AddOp = Plus | Minus
data MulOp = Times | Divide

tokens      :: String → [Token]
tokens []   = [EOF]
tokens ('+' : cs) = ADDOP Plus : tokens cs
tokens ('-' : cs) = ADDOP Minus : tokens cs
tokens ('*' : cs) = MULOP Times : tokens cs
tokens ('/' : cs) = MULOP Divide : tokens cs
tokens '(' : cs) = LPAREN : tokens cs
tokens ')' : cs) = RPAREN : tokens cs
tokens (c : cs)
  | isAlpha c = IDENT (c : n) : tokens cs'
  | otherwise = tokens cs
where (n, cs') = span isAlphaNum cs

```

Figure 2. Tokens for the expression grammar.

```

data Expr = RuleE→E+T Expr AddOp Term
           | RuleE→T Term
data Term = RuleT→T*F Term MulOp Factor
           | RuleT→F Factor
data Factor = RuleF→(E) Expr
             | RuleF→id String

```

Figure 3. Parse tree types for the expression grammar.

symbols may be computed by a fold over the parse tree. This fold may be automatically fused with the parser by a routine application of deforestation (Wad88), replacing the constructors  $\text{Rule}_{A \rightarrow \alpha}$  with functions that compute the values of the attributes (see also Section 6.4).

```

class ParseTree a where
  flatten :: a → [Token]
instance ParseTree Expr where
  flatten (RuleE→E+T e op t) = flatten e ++ [ADDOP op] ++ flatten t
  flatten (RuleE→T t)       = flatten t
instance ParseTree Term where
  flatten (RuleT→T*F t op f) = flatten t ++ [MULOP op] ++ flatten f
  flatten (RuleT→F f)       = flatten f
instance ParseTree Factor where
  flatten (RuleF→(E) e)      = [LPAREN] ++ flatten e ++ [RPAREN]
  flatten (RuleF→id n)      = [IDENT n]

```

Figure 4. Flattening parse trees for the expression grammar.

As a special case of this fold, we can define a family of functions  $\text{flatten}_\alpha$  for each string  $\alpha$ , with which we can recover the original sequence of tokens from a parse tree:

```

flattenα :: ᾱ → [Token]
flattenS† v = flattenS§ v
flattena va = [Ta va]
flattenA (RuleA→α vα) = flattenα vα
flattenε = []
flattenαβ vα vβ = flattenα vα ++ flattenβ vβ

```

This family of functions, together with the parse tree types, thus encode the grammar in Haskell. The initiality of parse trees captures the notion that the language defined is the least fixed point of the grammar in the subset ordering. Suitable definitions for the expression grammar are given in Figure 4. Note that the family of functions is represented by an overloaded function in Haskell.

### 3.4. NON-DETERMINISTIC PARSING

Parsing can now be described as an inverse of the function  $\text{flatten}_{S\dagger}$ . Of course, this function is not necessarily surjective or injective: some token strings may not be the flattening of any parse tree, while others may correspond to more than one parse tree, if the grammar is ambiguous. Although for practical applications we are primarily interested in unambiguous grammars, specifically LALR ones, it will be cleaner to develop our parsers for the general case. In general, a token string may be the flattening of a number of parse trees (possibly zero), so we define

```

type Parser r = [Token] → ℙ r

```



*Specification of parse* The parse function for a start symbol  $S$  will have type

$$\text{parse}_S :: \text{Parser } \bar{S}$$

We shall expect it to satisfy

$$\text{parse}_S \text{ ts} = \{ v_S \mid \text{flatten}_{S^\dagger} v_S \# \text{tr} = \text{ts} \} \quad (1)$$

Following the conventions discussed in Section 2, the variables  $v_S$  and  $\text{tr}$  are implicitly bound in the set comprehension. In this case, the suffix  $\text{tr}$  is actually immaterial, since  $\text{flatten}_{S^\dagger} v_S$  will end with the special token  $T_\S$  (EOF in the expression parser). Anything after that token will be ignored by the parser.

#### 4. Non-deterministic recursive descent parsing

We cannot expect a recursive definition of  $\text{parse}_S$ . Rather we shall define it in terms of a more general auxiliary function that parses a suffix of a production. Following convention we use items of the form  $A \rightarrow \alpha \cdot \beta$  to indicate this position. In the sequel we specify a parser for each item  $A \rightarrow \alpha \cdot \beta$ , called  $\text{state}_{A \rightarrow \alpha \cdot \beta}$ , with which we can define  $\text{parse}_S$ . We derive equations that this family of functions satisfies, and show that it is the least such family. We thus obtain a recursive characterization, which corresponds to recursive descent parsing.

*Specification of state* Suppose that during the parse we are at a point described by an item  $A \rightarrow \alpha \cdot \beta$ . We expect the remaining input to be parsable as a  $\beta$ , followed by further tokens. Hence the function to parse from this point is passed a parsing continuation, to be used after the  $\beta$  is seen. In parsing the string  $\beta$ , our parser will obtain attribute values of type  $\bar{\beta}$ , and these will be passed to the continuation. Hence the parsing function corresponding to  $A \rightarrow \alpha \cdot \beta$  has the signature

$$\text{state}_{A \rightarrow \alpha \cdot \beta} :: (\bar{\beta} \rightarrow \text{Parser } r) \rightarrow \text{Parser } r$$

and we shall require that it satisfies

$$\text{state}_{A \rightarrow \alpha \cdot \beta} \text{ k ts} = \bigcup \{ \text{k } v_\beta \text{ tr} \mid \text{flatten}_\beta v_\beta \# \text{tr} = \text{ts} \} \quad (2)$$

Note that  $\text{state}_{A \rightarrow \alpha \cdot \beta}$  is polymorphic in the result type  $r$ , which allows us to support parsers with multiple entry points (corresponding to grammars with multiple start symbols). Note further that  $\text{state}_{A \rightarrow \alpha \cdot \beta}$  neither depends on  $A$  nor on  $\alpha$ . Consequently, the  $\text{state}$  functions are

identical if the strings to the right of the dot are identical, that is,  $\text{state}_{A \rightarrow \alpha \cdot \beta} = \text{state}_{A' \rightarrow \alpha' \cdot \beta}$ . We return to this issue in Section 6.3.

*Derivation of parse* Parsing begins with the item  $S^\dagger \rightarrow \cdot S\$$ :

$$\begin{aligned}
& \text{parse}_S \text{ ts} \\
= & \langle \text{specification of parse} \rangle \\
& \{ v_S \mid \text{flatten}_{S^\dagger} v_S \# \text{tr} = \text{ts} \} \\
= & \langle \text{definition of flatten} \rangle \\
& \{ v_S \mid \text{flatten}_{S\$} v_S \# \text{tr} = \text{ts} \} \\
= & \langle \text{sets} \rangle \\
& \bigcup \{ \{ v_S \} \mid \text{flatten}_{S\$} v_S \# \text{tr} = \text{ts} \} \\
= & \langle \text{introduce accept—see below} \rangle \\
& \bigcup \{ \text{accept } v_S \text{ tr} \mid \text{flatten}_{S\$} v_S \# \text{tr} = \text{ts} \} \\
= & \langle \text{specification of state and } \overline{S\$} = \overline{S} \rangle \\
& \text{state}_{S^\dagger \rightarrow \cdot S\$} \text{ accept ts}
\end{aligned}$$

where the function `accept` is defined as

$$\begin{aligned}
\text{accept} & \quad :: r \rightarrow \text{Parser } r \\
\text{accept } v \text{ tr} & = \{ v \}
\end{aligned}$$

and so we have

$$\text{parse}_S = \text{state}_{S^\dagger \rightarrow \cdot S\$} \text{ accept}$$

*Derivation of state* We proceed by case analysis of the string to the right of the dot.

**Case  $A \rightarrow \alpha \cdot$ :** Consider first items with an empty string to the right of the dot. In this case, we can use the parsing continuation to parse the rest of `ts`:

$$\begin{aligned}
& \text{state}_{A \rightarrow \alpha \cdot} \text{ k ts} \\
= & \langle \text{specification of state} \rangle \\
& \bigcup \{ \text{k } v_\epsilon \text{ tr} \mid \text{flatten}_\epsilon v_\epsilon \# \text{tr} = \text{ts} \} \\
= & \langle \text{definition of flatten} \rangle \\
& \bigcup \{ \text{k ts} \} \\
= & \langle \text{sets} \rangle \\
& \text{k ts}
\end{aligned}$$

In a conventional setting this corresponds to a reduce step.

**Case  $A \rightarrow \alpha \cdot X\beta$ :** If there is at least one symbol to the right of the dot, we have

$$\begin{aligned}
& \text{state}_{A \rightarrow \alpha \cdot X\beta} \text{ k ts} \\
= & \langle \text{specification of state} \rangle \\
& \bigcup \{ \text{k } v_{X\beta} \text{ tr} \mid \text{flatten}_{X\beta} v_{X\beta} \# \text{tr} = \text{ts} \} \\
= & \langle \text{definition of flatten} \rangle \\
& \bigcup \{ \text{k } v_X v_\beta \text{ tr} \mid \text{flatten}_X v_X \# \text{flatten}_\beta v_\beta \# \text{tr} = \text{ts} \} \\
= & \langle \text{sets} \rangle \\
& \bigcup \{ \bigcup \{ \text{k } v_X v_\beta \text{ tr} \mid \text{flatten}_\beta v_\beta \# \text{tr} = \text{tr}' \} \\
& \quad \mid \text{flatten}_X v_X \# \text{tr}' = \text{ts} \} \\
= & \langle \text{specification of state} \rangle \\
& \bigcup \{ \text{state}_{A \rightarrow \alpha X \cdot \beta} (\text{k } v_X) \text{ tr}' \mid \text{flatten}_X v_X \# \text{tr}' = \text{ts} \} \\
= & \langle \text{introduce goto—see below} \rangle \\
& \bigcup \{ \text{goto}_{A \rightarrow \alpha \cdot X\beta} \text{ k } v_X \text{ tr}' \mid \text{flatten}_X v_X \# \text{tr}' = \text{ts} \} \\
= & \langle \text{introduce shift—see below} \rangle \\
& \text{shift}_X (\text{goto}_{A \rightarrow \alpha \cdot X\beta} \text{ k}) \text{ ts}
\end{aligned}$$

where functions  $\text{goto}_{A \rightarrow \alpha \cdot X\beta}$  and  $\text{shift}_X$  are specified by

$$\begin{aligned}
\text{goto}_{A \rightarrow \alpha \cdot X\beta} & \quad :: (\overline{X\beta} \rightarrow \text{Parser } r) \rightarrow (\overline{X} \rightarrow \text{Parser } r) \\
\text{goto}_{A \rightarrow \alpha \cdot X\beta} \text{ k } v_X & = \text{state}_{A \rightarrow \alpha X \cdot \beta} (\text{k } v_X) \\
\text{shift}_X & \quad :: (\overline{X} \rightarrow \text{Parser } r) \rightarrow \text{Parser } r \\
\text{shift}_X \text{ g ts} & = \bigcup \{ \text{g } v_X \text{ tr} \mid \text{flatten}_X v_X \# \text{tr} = \text{ts} \}
\end{aligned}$$

As further progress depends on the nature of  $X$ , we have introduced  $\text{shift}$  as a name for the branch point. We also introduce  $\text{goto}$  to cast  $\text{shift}$  in a more convenient form for further manipulation:  $\text{goto}_{A \rightarrow \alpha \cdot X\beta}$  feeds attributes to the continuation and enters the state after  $X$ .

*Derivation of shift* We proceed by case analysis on  $X$ .

**Case  $X = c$ :** When  $X$  is a terminal symbol  $c$ , we have

$$\begin{aligned}
& \text{shift}_c \text{ g ts} \\
= & \langle \text{specification of shift} \rangle \\
& \bigcup \{ \text{g } v_c \text{ tr} \mid \text{flatten}_c v_c \# \text{tr} = \text{ts} \} \\
= & \langle \text{definition of flatten} \rangle \\
& \bigcup \{ \text{g } v_c \text{ tr} \mid [\text{T}_c v_c] \# \text{tr} = \text{ts} \} \\
= & \langle \text{definition of ‘\#’} \rangle \\
& \bigcup \{ \text{g } v_c \text{ tr} \mid \text{T}_c v_c : \text{tr} = \text{ts} \} \\
= & \langle \text{case expression} \rangle \\
& \text{case ts of } (\text{T}_c v_c : \text{tr}) \rightarrow \text{g } v_c \text{ tr} \\
& \quad \quad \quad \_ \quad \quad \rightarrow \emptyset
\end{aligned}$$

In a conventional setting this corresponds to a shift step.

**Case  $X = C$ :** When  $X$  is a nonterminal symbol  $C$ , we have

$$\begin{aligned}
& \text{shift}_C \text{ g ts} \\
= & \langle \text{specification of shift} \rangle \\
& \bigcup \{ \text{g } v_C \text{ tr} \mid \text{flatten}_C v_C \# \text{tr} = \text{ts} \} \\
= & \langle \text{definition of flatten} \rangle \\
& \bigcup \{ \text{g } (\text{Rule}_{C \rightarrow \gamma} v_\gamma) \text{ tr} \mid C \rightarrow \gamma \in \mathbf{P} \wedge \text{flatten}_\gamma v_\gamma \# \text{tr} = \text{ts} \} \\
= & \langle \text{introduce reduce—see below} \rangle \\
& \bigcup \{ \text{reduce}_{C \rightarrow \gamma} \text{ g } v_\gamma \text{ tr} \mid C \rightarrow \gamma \in \mathbf{P} \wedge \text{flatten}_\gamma v_\gamma \# \text{tr} = \text{ts} \} \\
= & \langle \text{sets} \rangle \\
& \bigcup \{ \bigcup \{ \text{reduce}_{C \rightarrow \gamma} \text{ g } v_\gamma \text{ tr} \mid \text{flatten}_\gamma v_\gamma \# \text{tr} = \text{ts} \} \\
& \quad \mid C \rightarrow \gamma \in \mathbf{P} \} \\
= & \langle \text{specification of state} \rangle \\
& \bigcup \{ \text{state}_{C \rightarrow \cdot \gamma} (\text{reduce}_{C \rightarrow \gamma} \text{ g}) \text{ ts} \mid C \rightarrow \gamma \in \mathbf{P} \}
\end{aligned}$$

where we define  $\text{reduce}_{C \rightarrow \gamma}$  as the continuation-passing counterpart of  $\text{Rule}_{C \rightarrow \gamma}$ :

$$\begin{aligned}
\text{reduce}_{C \rightarrow \gamma} & \quad :: (\bar{C} \rightarrow \text{Parser } r) \rightarrow (\bar{\gamma} \rightarrow \text{Parser } r) \\
\text{reduce}_{C \rightarrow \gamma} \text{ g } v_\gamma & = \text{g } (\text{Rule}_{C \rightarrow \gamma} v_\gamma)
\end{aligned}$$

This step, the only remaining source of nondeterminism, has no counterpart in a conventional setting. We shall transform it further in the next section. For the expression parser, the **reduce** functions are given in Figure 4. Note that each constructor  $\text{Rule}_{C \rightarrow \gamma}$  is mentioned once only. As remarked in Section 3.3, when the parser is fused with a fold over a parse tree and deforested, these constructors will be replaced by actions.

To summarize, Figure 6 lists the derived properties of **state** and of the auxiliary functions we have introduced.

In general, **state** will not be the sole family of functions satisfying these equations, but we can show that it is the least such family under the set inclusion ordering, lifted pointwise to functions. Let  $\text{state}'$ ,  $\text{shift}'$  and  $\text{goto}'$  be another such family. Then we wish to establish that  $\text{state} \subseteq \text{state}'$  and  $\text{shift} \subseteq \text{shift}'$ , or equivalently

$$k v_\beta \text{ tr} \subseteq \text{state}'_{A \rightarrow \alpha \cdot \beta} k (\text{flatten}_\beta v_\beta \# \text{tr}) \quad (3)$$

$$\text{g } v_X \text{ ts} \subseteq \text{shift}'_X \text{ g } (\text{flatten}_X v_X \# \text{ts}) \quad (4)$$

The first set inclusion states that  $\text{state}'$  produces every parse tree by parsing its flattening and possibly others. The proof is similar in structure to the derivation: We show (3) by induction over the depth of the parse forest  $v_\beta$  with (4) serving as an intermediate result. The details of the proof are omitted.

If the grammar is not left-recursive, the equations of Figure 6 are, in fact, executable, acting as a nondeterministic LL(0) parser. The parser



## 5. Non-deterministic recursive ascent parsing

As noted above, the equations for items of the forms  $A \rightarrow \alpha \cdot$  and  $A \rightarrow \alpha \cdot c\beta$  are directly implementable, corresponding to the reduce and shift actions of a traditional LR parser. It remains to eliminate actions of the form  $A \rightarrow \alpha \cdot C\beta$  by expanding them to reduce or shift actions. For this step, we operate on the continuation-passing functions `state` and `shiftC` expressed in a point-free style:

$$\begin{aligned} \text{state}_{A \rightarrow \alpha \cdot} &= \text{id} \\ \text{state}_{A \rightarrow \alpha \cdot X\beta} &= \text{shift}_X \circ \text{goto}_{A \rightarrow \alpha \cdot X\beta} \\ \text{shift}_C &= \bigcup \{ \text{state}_{C \rightarrow \cdot \gamma} \circ \text{reduce}_{C \rightarrow \gamma} \mid C \rightarrow \gamma \in \mathbf{P} \} \end{aligned}$$

where in the last equation union is lifted pointwise to function spaces:  $(f \cup g) x = f x \cup g x$ . Each `state` and each `shiftC` expands to the union of a possibly infinite set of compositions of finite sequences of `goto`, `shiftc` and `reduce` functions. We shall denote these sets of sequences by nonterminal symbols `State` and `ShiftC` respectively, defined by a left-linear grammar corresponding to the above equations:

$$\begin{aligned} \text{State}_{A \rightarrow \alpha \cdot} &\rightarrow \epsilon \\ \text{State}_{A \rightarrow \alpha \cdot X\beta} &\rightarrow \text{Shift}_X \text{ goto}_{A \rightarrow \alpha \cdot X\beta} \\ \text{Shift}_c &\rightarrow \text{shift}_c \\ \text{Shift}_C &\rightarrow \text{State}_{C \rightarrow \cdot \gamma} \text{ reduce}_{C \rightarrow \gamma} \quad C \rightarrow \gamma \in \mathbf{P} \end{aligned}$$

From this perspective, the task of eliminating left recursion reduces to the standard procedure of transforming this left-linear grammar into an equivalent right-linear grammar (see Section 2).

There is, however, a caveat. The transformation only preserves the meaning of the start symbol. Therefore, we confine our attention to `Stateq` for an arbitrary, but fixed kernel item  $q$ . The derivation below must then be repeated for every kernel item. (We need not consider non-kernel items though.) We furthermore restrict the grammar to symbols reachable from `Stateq` as this does not change the language denoted by `Stateq`. Then `Stateq` is given by the following productions (for readability we have added a new start symbol `Start`):

$$\begin{aligned} \text{Start} &\rightarrow \text{State}_q \\ \text{State}_i &\rightarrow \epsilon && i = A \rightarrow \alpha \cdot \in q^* \\ \text{State}_i &\rightarrow \text{Shift}_X \text{ goto}_i && i \in q^*/X \\ \text{Shift}_c &\rightarrow \text{shift}_c && c \in \text{next}(q^*) \\ \text{Shift}_C &\rightarrow \text{State}_{C \rightarrow \cdot \gamma} \text{ reduce}_{C \rightarrow \gamma} && C \rightarrow \cdot \gamma \in q^+ \end{aligned}$$

The equivalent right-linear grammar is then

$$\begin{array}{lll}
\overline{\text{State}}_q & \rightarrow \epsilon & \\
\overline{\text{Start}} & \rightarrow \overline{\text{State}}_i & i = A \rightarrow \alpha \cdot \in q^* \\
\overline{\text{Shift}}_X & \rightarrow \text{goto}_i \overline{\text{State}}_i & i \in q^*/X \\
\overline{\text{Start}} & \rightarrow \text{shift}_c \overline{\text{Shift}}_c & c \in \text{next}(q^*) \\
\overline{\text{State}}_{C \rightarrow \cdot \gamma} & \rightarrow \text{reduce}_{C \rightarrow \gamma} \overline{\text{Shift}}_C & C \rightarrow \cdot \gamma \in q^+
\end{array}$$

Reverting the right-linear grammar back into a program we obtain ( $\overline{\text{state}}_q$  corresponds to  $\overline{\text{Start}}$ ):

$$\begin{array}{ll}
\overline{\text{state}}_q & = \bigcup \{ \text{shift}_c \circ \overline{\text{shift}}_c \mid c \in \text{next}(q^*) \} \cup \\
& \quad \bigcup \{ \overline{\text{state}}_i \mid i = A \rightarrow \alpha \cdot \in q^* \} \\
\overline{\text{state}}_q & = \text{id} \\
\overline{\text{state}}_{C \rightarrow \cdot \gamma \in q^+} & = \text{reduce}_{C \rightarrow \gamma} \circ \overline{\text{shift}}_C \\
\overline{\text{shift}}_X & = \bigcup \{ \text{goto}_i \circ \overline{\text{state}}_i \mid i \in q^*/X \}
\end{array}$$

Note that the transformation relies on the continuity of  $\text{shift}$ ,  $\text{reduce}$  and  $\text{goto}$ . Furthermore, it is important to note that the two equations for  $\overline{\text{state}}$  do not overlap because  $q \notin q^+$  for any kernel item  $q$ .

Now, for the following it is more convenient to rewrite the definition above into an applicative form.

$$\begin{array}{ll}
\overline{\text{state}}_q k & = \bigcup \{ \text{shift}_c (\overline{\text{shift}}_c k) \mid c \in \text{next}(q^*) \} \cup \\
& \quad \bigcup \{ \overline{\text{state}}_i k \mid i = A \rightarrow \alpha \cdot \in q^* \} \\
\overline{\text{state}}_q k & = k \\
\overline{\text{state}}_{C \rightarrow \cdot \gamma \in q^+} k & = \text{reduce}_{C \rightarrow \gamma} (\overline{\text{shift}}_C k) \\
\overline{\text{shift}}_X k & = \bigcup \{ \text{goto}_i (\overline{\text{state}}_i k) \mid i \in q^*/X \}
\end{array}$$

Since we are aiming at a single definition for  $\overline{\text{state}}_q$ , we turn  $\overline{\text{state}}$  and  $\overline{\text{shift}}$  into local definitions setting  $k_i = \overline{\text{state}}_i k$  and  $g_X = \overline{\text{shift}}_X k$ :

$$\begin{array}{ll}
\overline{\text{state}}_q k & = \bigcup \{ \text{shift}_c g_c \mid c \in \text{next}(q^*) \} \cup \\
& \quad \bigcup \{ k_i \mid i = A \rightarrow \alpha \cdot \in q^* \} \\
\text{where } k_q & = k \\
k_{C \rightarrow \cdot \gamma \in q^+} & = \text{reduce}_{C \rightarrow \gamma} g_C \\
g_X & = \bigcup \{ \text{goto}_i k_i \mid i \in q^*/X \}
\end{array}$$

The last transformation is also known as  $\lambda$ -dropping (Dan99). We can finally expand  $g_X$  by calculating

$$\begin{array}{l}
g_X \nu_X \\
= \langle \text{definition of } g_X \rangle \\
\bigcup \{ \text{goto}_i k_i \nu_X \mid i = A \rightarrow \alpha \cdot X\beta \in q^* \} \\
= \langle \text{definition of } \text{goto}_i \rangle \\
\bigcup \{ \text{state}_{A \rightarrow \alpha X \cdot \beta} (k_i \nu_X) \mid i = A \rightarrow \alpha \cdot X\beta \in q^* \} \\
= \langle \text{definition of } Q+X \text{ and } Q/X \rangle \\
\bigcup \{ \text{state}_{i+X} (k_i \nu_X) \mid i \in q^*/X \}
\end{array}$$

$$\begin{aligned}
\text{state}_{A \rightarrow \alpha \cdot \beta} &:: (\bar{\beta} \rightarrow \text{Parser } r) \rightarrow \text{Parser } r \\
\text{state}_q \text{ k}_q \text{ ts} &= \bigcup \{ \text{shift}_c \text{ g}_c \text{ ts} \mid c \in \text{next}(q^*) \} \cup \\
&\quad \bigcup \{ \text{k}_i \text{ ts} \mid i = A \rightarrow \alpha \cdot \in q^* \} \\
\text{where } \text{k}_{C \rightarrow \cdot \gamma \in q^+} &= \text{reduce}_{C \rightarrow \gamma} \text{ g}_C \\
\text{g}_X \text{ v}_X &= \bigcup \{ \text{state}_{i+X} (\text{k}_i \text{ v}_X) \mid i \in q^*/X \}
\end{aligned}$$

Figure 7. Non-deterministic LR(0) parser for a kernel item  $q$ .

$$\begin{aligned}
\text{state}_Q &:: (\bar{\beta} \rightarrow \text{Parser } r \mid A \rightarrow \alpha \cdot \beta \in Q) \rightarrow \text{Parser } r \\
\text{state}_Q (\text{k}_q \mid q \in Q) \text{ ts} &= \bigcup \{ \text{shift}_c \text{ g}_c \text{ ts} \mid c \in \text{next}(Q^*) \} \cup \\
&\quad \bigcup \{ \text{k}_i \text{ ts} \mid i = A \rightarrow \alpha \cdot \in Q^* \} \\
\text{where } \text{k}_{C \rightarrow \cdot \gamma \in Q^+} &= \text{reduce}_{C \rightarrow \gamma} \text{ g}_C \\
\text{g}_X \text{ v}_X &= \text{state}_{Q^*+X} (\text{k}_i \text{ v}_X \mid i \in Q^*/X)
\end{aligned}$$

Figure 8. Non-deterministic LR(0) parser for a set of kernel items  $Q$ .

We have derived the parser of Figure 7. The parser for our running example, the expression grammar, can be found in Appendix B. Although our parsers are not yet deterministic, a relationship to traditional table-driven parsers can be noted: the right-hand side of  $\text{state}_q$  corresponds to an entry in the action table, the local function  $\text{g}_X$  to an entry of the goto table.

## 6. Deterministic recursive ascent parsing

### 6.1. PURSUING SETS OF ITEMS IN PARALLEL

The first step towards a deterministic parser is to collect sets of parsers for individual kernel items into a parser for a set  $Q$  of kernel items, with a corresponding tuple of parsing continuations:

$$\text{state}_Q (\text{k}_q \mid q \in Q) \text{ ts} = \bigcup \{ \text{state}_q \text{ k}_q \text{ ts} \mid q \in Q \}$$

This corresponds to the classical subset construction for converting a non-deterministic finite automaton to a deterministic one. Since everything distributes over union, we obtain the parser of Figure 8.

Using the properties of **case**, the unions of  $\text{shift}_c$  functions may furthermore be combined into a single **case** statement:

$$\begin{aligned}
\text{shift}_{c_1} \text{ g}_1 \text{ ts} \cup \dots \cup \text{shift}_{c_n} \text{ g}_n \text{ ts} &= \text{case ts of} \\
&\quad \text{T}_{c_1} \text{ v}_{c_1} : \text{tr} \rightarrow \text{g}_1 \text{ v}_{c_1} \text{ tr} \\
&\quad \dots \\
&\quad \text{T}_{c_n} \text{ v}_{c_n} : \text{tr} \rightarrow \text{g}_n \text{ v}_{c_n} \text{ tr} \\
&\quad - \quad \quad \quad \rightarrow \emptyset
\end{aligned}$$



The state functions in our parsers take parsing continuations as arguments, one for each kernel item in the corresponding state of the LR automaton. Traditional parsers also maintain a stack of attribute values, popping the required number of values when a reduction occurs. In contrast, in our parsers the parsing continuations are applied to these attributes as they become available. In effect, we have a short value stack for each continuation. This seems expensive, but the number of continuations of a state is equal to the number of kernel items, and in typical parsers this is small.

## 6.2. USING LOOK-AHEAD

The remaining non-determinism of the parser in Figure 8 arises when a state contains both shifts and reductions (shift-reduce conflict) or more than one reduction (reduce-reduce conflict). By definition, these conflicts do not occur if the grammar is LR(0). For these grammars we have constructed a deterministic parser: the set of results is either empty or has a single element. We can represent such sets using Haskell's `Maybe` type. For non-LR(0) grammars we can use look-ahead information to resolve these conflicts. For example, in the expression parser the item set  $\{E \rightarrow E+T \cdot, T \rightarrow T \cdot *F\}$  yields the state function

```
state6 :: Parser r → (MulOp → Factor → Parser r) → Parser r
state6 kE→E+T. kT→T.*F ts =
  (case ts of MULOP v : tr → state8 (kT→T.*F v) tr
    _                → ∅)
  ∪ kE→E+T. ts -- shift-reduce conflict
```

The look-ahead set for the reduce item  $k_{E \rightarrow E+T}$  is  $\{\text{RPAREN}, \text{ADDOP}, \text{EOF}\}$ . Using this information we can refine the case expression to

```
state6 :: Parser r → (MulOp → Factor → Parser r) → Parser r
state6 kE→E+T. kT→T.*F ts =
  case ts of MULOP v : tr → state8 (kT→T.*F v) tr
            ADDOP v : _  → kE→E+T. ts
            RPAREN  : _  → kE→E+T. ts
            EOF     : _  → kE→E+T. ts
            _       : _  → ∅
```

Eliminating reductions in this way removes many conflicts. If the grammar is LALR(1), it makes the parser deterministic. Alternatively, we can perform the reduction on all symbols other than `MULOP`.

```
state6 :: Parser r → (MulOp → Factor → Parser r) → Parser r
state6 kE→E+T. kT→T.*F ts =
  case ts of MULOP v : tr → state8 (kT→T.*F v) tr
            _           → kE→E+T. ts
```

This choice corresponds to a standard optimization of table-based LR parsers (ASU86), where the most common reduce action is generalized. The resulting expression parser is given in Appendix C.

### 6.3. SMALLER PARSERS

For realistic grammars, the parsers generated by the above method are fast, but quite large, possibly stretching the compiler or the available memory. However, several simple steps can drastically reduce the size of the parser. Since our parsers eschew tables, exposing the control flow, a reasonable optimizing compiler would perform many of these steps anyway. However, our concern is that the input program would be too large for the compiler, so they should be done by the parser generator.

Firstly, because in our parsers the parsing continuations are already applied to past values, the only part of a kernel item  $A \rightarrow \alpha \cdot \beta$  that matters is  $\beta$ . If these parts of the items, and the associated look-aheads, are identical, the states can be merged. For example, in the parser of Appendix C, states 6 and 7 are identical. The whole items, however, do not match, so the states could not be merged in a table-based parser.

The extreme case, where a state consists of a single item, and that of the form  $A \rightarrow \alpha \cdot$ , is quite common. In this case, the state is an identity function, and thus may be eliminated. In the example parser, states 3, 9, 10, 12 and 13 may be eliminated in this way.

Now note that states 1, 5, 8 and 11 have identical actions, differing only in the goto functions they use. So one could abstract out the common action part, and parameterize it by the goto functions and continuations used. This is a common situation, and corresponds to a standard optimization for table-based parsers (ASU86). In this example, this common part is equivalent to `state8` (after the previous optimization):

$$\begin{aligned} \text{state}_8 &:: (\text{Factor} \rightarrow \text{Parser } r) \rightarrow \text{Parser } r \\ \text{state}_8 \text{ } k_{T \rightarrow T^* \cdot F} \text{ } ts &= \\ \text{case } ts \text{ of LPAREN } : \text{tr} &\rightarrow \text{state}_{11} (\text{reduce}_{F \rightarrow (E)} g_F) \text{tr} \\ \text{IDENT } v : \text{tr} &\rightarrow \text{reduce}_{F \rightarrow \text{id}} k_{T \rightarrow T^* \cdot F} v \text{tr} \\ - &\rightarrow \emptyset \end{aligned}$$

We can go further: the definition of  $g_T$  and  $g_F$  in states 1 and 11 is the same as in state 5, so they may be defined in terms of it. (This could also be done in a table-based parser, but would involve slowing down the selection of the goto function.) The resulting expression parser is given in Appendix D.

```

type Environment = [(String, Int)]
type Expr'       = Environment → Int
type Term'      = Environment → Int
type Factor'    = Environment → Int
ruleE→E+T      :: Expr' → AddOp → Term' → Expr'
ruleE→E+T e Plus t env = e env + t env
ruleE→E+T e Minus t env = e env - t env
ruleE→T        :: Term' → Expr'
ruleE→T        = id
ruleT→T*F      :: Term' → MulOp → Factor' → Term'
ruleT→T*F t Times f env = t env * f env
ruleT→T*F t Divide f env = t env `div` f env
ruleT→F       :: Factor' → Term'
ruleT→F       = id
ruleF→(E)     :: Expr' → Factor'
ruleF→(E)     = id
ruleF→id      :: String → Factor'
ruleF→id n env = fromJust (lookup n env)

```

Figure 9. Attribute functions for an environment-based evaluator.

#### 6.4. ATTRIBUTES

In many applications, we will want to apply a fold to the parse tree returned by our parser. A simple form of deforestation (GLP83) is to redefine the attribute type  $\bar{A}$  of each nonterminal  $A$  and replace each  $\text{Rule}_{A \rightarrow \alpha}$  constructor with a corresponding function  $\text{rule}_{A \rightarrow \alpha}$ . Since the  $\text{Rule}$  constructors occur only in the  $\text{reduce}$  functions, only these need be changed, in addition to redefining the types.

For example, if we wished to evaluate a parsed expression in an environment mapping identifiers to integer values, we could use the definitions of Figure 9.

### 7. Related work on functional LR parsers

Leermakers, Augsteijn and Kruseman Aretz (LAK92; Lee93; Aug93) defined several functional parsers, including an LR parser consisting of functions  $[Q]$  for a set  $Q$  of kernel items specified as

$$[Q] w = \{ (A \rightarrow \alpha \cdot \beta, w_2) \mid A \rightarrow \alpha \cdot \beta \in Q \wedge \beta \rightarrow^* w_1 \wedge w = w_1 w_2 \}$$

Many nondeterministic recursive functions, including the functions  $[Q]$  above, are amenable to a group of general transformations studied

$$\begin{aligned}
[Q] &:: [\text{Token}] \rightarrow \mathbb{P}(\text{Item}, [\text{Token}]) \\
[Q] w &= \overline{[Q]}_c w' \cup \\
&\quad \{ (A \rightarrow \alpha \cdot, w) \mid A \rightarrow \alpha \cdot \in Q \} \cup \\
&\quad \cup \{ \overline{[Q]}_C w \mid C \rightarrow \cdot \in Q^+ \} \\
&\quad \mathbf{where} (c : w') = w \\
\overline{[Q]}_X &:: [\text{Token}] \rightarrow \mathbb{P}(\text{Item}, [\text{Token}]) \\
\overline{[Q]}_X w &= \{ (i, w') \mid i \in Q/X \wedge (i+X, w') \in [Q^*+X] w \} \cup \\
&\quad \cup \{ \overline{[Q]}_C w' \mid j \in Q^+/X \wedge (j+X, w') \in [Q^*+X] w \}
\end{aligned}$$

Figure 10. Direct-style functional parser of Leermakers.

by Augusteijn (Aug93). These transformations lead to the definition of Figure 10, where we have omitted handling of attributes to avoid clutter. Note a minor complication in the definition of  $[Q]$ , where kernel reduce items  $A \rightarrow \alpha \cdot$  and predict reduce items  $C \rightarrow \cdot$  are treated differently.

This version was further refined to an imperative implementation, providing a derivation of the parsers of Kruseman Aretz (Kru88) and Pennello (Pen86).

Leermakers also defined a continuation-passing version, where the continuation is passed the pair  $(A \rightarrow \alpha \cdot \beta, w_2)$ , and then winds back through  $|\alpha|$  continuations and dispatches on the nonterminal  $A$ , corresponding to the traditional reduction and goto function.

Sperber and Thiemann (ST00) also derived a continuation-passing version, by applying partial evaluation techniques to the parser of Figure 10. Their variant passes a list of continuations, corresponding to the stack of states (goto functions) in traditional parsers. They realized that only a finite number of these were needed (the largest  $|\alpha|$  for any item  $A \rightarrow \alpha \cdot \beta$  in  $Q$ ) so they could be passed as separate arguments. A consequence of these transformations is that kernel and predict items may be handled uniformly, as in our parser. They handled attributes using a stack of values, which may be similarly truncated. Sperber and Thiemann worked in an untyped setting, but their parser may be recast in typed form as in Figure 11. Each  $k_i$  here corresponds to a row of the goto table of a traditional parser, which is usually quite sparse. On reduction, the parser selects the continuation, which then dispatches over the nonterminal. In comparison, our parser passes entries of the goto table as separate arguments.

It is possible to remove the interpretive overhead from Sperber and Thiemann's parser by fusing  $k_{|\gamma|} \circ N_C$ . In this way, we obtain the parser

$$\begin{aligned}
\mathbf{data} \text{ NonTerminal} &= \sum_{A \in \mathbf{N}^\dagger} \mathbf{N}_A \bar{A} \\
[Q] &:: (\text{NonTerminal} \rightarrow \text{Parser } r)^{|\alpha|} \rightarrow \bar{\alpha} \rightarrow \text{Parser } r \\
[Q] \ k_1 \ \dots \ k_{|\alpha|} \ v_\alpha \ \mathbf{ts} &= \bigcup \{ \text{shift}_c \ g_c \ \mathbf{ts} \mid c \in \text{next}(Q^*) \} \\
&\quad \cup \bigcup \{ \text{reduce}_{C \rightarrow \gamma} (k_{|\gamma|} \circ \mathbf{N}_C) \ v_\gamma \ \mathbf{ts} \mid C \rightarrow \gamma \cdot \in Q^* \} \\
\mathbf{where} \ g_X \ v_X &= [Q^* + X] \ k_0 \ k_1 \ \dots \ k_{|\alpha'|} \ v_{\alpha'} \ v_X \\
k_0 (\mathbf{N}_A \ v_A) &= g_A \ v_A \quad A \in \text{next}(Q^*) \\
\alpha &\text{ maximal such that } A \rightarrow \alpha \cdot \beta \in Q \\
\alpha' &\text{ maximal such that } A' \rightarrow \alpha' \cdot X \beta' \in Q^*
\end{aligned}$$

Figure 11. Typed version of Sperber and Thiemann's parser for a set of kernel items  $Q$ .

$$\begin{aligned}
[Q] &:: (\bar{A} \rightarrow \text{Parser } r \mid A \rightarrow \alpha \cdot \beta \in Q) \rightarrow \bar{\alpha}_{\max} \rightarrow \text{Parser } r \\
[Q] \ (k_i \mid i \in Q) \ v_{\alpha_{\max}} \ \mathbf{ts} &= \bigcup \{ \text{shift}_c \ g_c \ \mathbf{ts} \mid c \in \text{next}(Q^*) \} \\
&\quad \cup \bigcup \{ \text{reduce}_{C \rightarrow \gamma} \ k_{C \rightarrow \gamma} \ v_\gamma \ \mathbf{ts} \mid C \rightarrow \gamma \cdot \in Q^* \} \\
\mathbf{where} \ g_X \ v_X &= [Q^* + X] \ (k_i \mid i \in Q^*/X) \ v_{\alpha'_{\max}} \ v_X \\
k_{C \rightarrow \cdot \gamma} &= g_C \\
\alpha_{\max} &\text{ maximal such that } A \rightarrow \alpha \cdot \beta \in Q \\
\alpha'_{\max} &\text{ maximal such that } A' \rightarrow \alpha' \cdot X \beta' \in Q^*
\end{aligned}$$

Figure 12. Intermediate parser for a set of kernel items  $Q$ .

of Figure 12. From this, we could then obtain the parser of Figure 8 by applying to continuations to the values as they become available, instead of passing the top portion of the stack.

Another scheme for typing LR parsers was presented by Pottier and Régis-Gianas (PRG05). They also address the interpretive overhead of the union type for nonterminals, but their solution requires a type system extended with inductive type families (also known as generalized algebraic data types).

A different approach is pursued by Pepper (Pep04). Though recursive descent formulations of grammars are not always directly executable, Pepper shows that if the underlying language is  $\text{LR}(k)$ , the programs may be transformed into equivalent recursive descent parsers, extended with a delayed choice operator.

## 8. Measurements

The purpose of this section is to compare the time and space behaviour of our parsers with that of traditional table-based ones. To get

grammar <sup>a</sup>	parser	src <sup>b</sup>	bin <sup>c</sup>	prg <sup>d</sup>	time <sup>e</sup>	% <sup>f</sup>	% <sup>g</sup>
expr 1.0M 0.100s	<b>frown -cs</b>	5.0K	18K	482K	2.184s	101	101
	<b>frown</b>	5.7K	18K	482K	2.160s	100	100
	<b>frown -cc</b>	4.3K	20K	484K	2.264s	105	105
	<b>frown -cg</b>	5.5K	25K	491K	2.480s	115	116
	<b>happy</b>	12K	27K	492K	6.784s	314	324
	<b>happy -cg</b>	14K	19K	482K	4.380s	203	208
	<b>happy -acg</b>	14K	16K	488K	4.772s	221	227
Oberon 2.9M 1.492s	<b>frown -cs</b>	114K	177K	1.2M	2.568s	100	100
	<b>frown</b>	169K	341K	1.4M	2.616s	102	104
	<b>frown -cc</b>	77K	242K	1.3M	2.620s	102	105
	<b>frown -cg</b>	118K	237K	1.3M	2.636s	103	106
	<b>happy</b>	87K	215K	1.2M	3.700s	144	205
	<b>happy -cg</b>	109K	122K	1.1M	4.180s	163	250
	<b>happy -acg</b>	84K	122K	1.1M	3.736s	145	209
Haskell 1.9M 3.016s	<b>frown -cs</b>	414K	1.2M	2.6M	4.172s	100	100
	<b>frown</b>	732K	2.1M	3.4M	4.272s	102	109
	<b>frown -cc</b>	298K	1.2M	2.6M	4.256s	102	107
	<b>frown -cg</b>	307K	1.0M	2.3M	4.204s	101	103
	<b>happy</b>	207K	666K	2.0M	5.068s	121	178
	<b>happy -cg</b>	243K	329K	1.7M	5.836s	140	244
	<b>happy -acg</b>	163K	205K	1.5M	5.512s	132	216

<sup>a</sup>grammar, size of the input to the parser, lexing time (minimum of 4 runs)

<sup>b</sup>size of the generated parser (source) <sup>c</sup>size of the generated parser (binary)

<sup>d</sup>size of the executable (including lexer and pretty printer)

<sup>e</sup>execution time (minimum of 4 runs, including lexing and pretty printing)

<sup>f</sup>time relative to the fastest parser (including lexing and pretty printing)

<sup>g</sup>time relative to the fastest parser (excluding lexing)

Figure 13. Benchmark results.

a somewhat broader picture we also include alternative non-table based parsers in the competition. Table 13 lists the results of the benchmarks, which were conducted on an AMD Athlon 64 3000+ with 2GB of main memory. We consider three grammars, a tiny, a medium-sized and a large grammar:

1. the expression grammar of Example 1 (7 productions, 13 states),
2. a grammar for Oberon (175 productions, 288 states),
3. a grammar for Haskell 98 (277 productions, 482 states).

The parsers for these grammars have been generated using two parser generators for Haskell: **happy** (GM05) is a yacc-like parser generator, **frown** (Hin05) is an LALR( $k$ ) parser generator developed by the first

author of the present paper. The generated parsers were compiled using the Glasgow Haskell Compiler, GHC 6.4.1 (The05), with optimizations, `-O`, turned on. Here are the salient features of the generated parsers:

1. `frown -cs` : the parsers developed in this paper;
2. `frown` : a table-free parser that uses an explicit stack of *state transitions* (represented by a tailor-made union type (DS00));
3. `frown -cc` : like 2 but with a more conventional stack type (a tailor-made union type for an alternation of states and attributes);
4. `frown -cg` : like 2 with the stack represented by a nested pair type (similar to the version of Figure 11 but passing the whole stack rather than a window of the top portion);
5. `happy` : a table-free parser that uses a polymorphic stack and a universal type for attributes (the union of all attribute types);
6. `happy -cg` : like 5, but using several GHC extensions including unboxed integers and unsafe type coercions (to avoid the need for a union type);
7. `happy -acg`: a table-based parser that uses the same GHC extensions as 6.

The different parsers for a grammar share the same lexical analyzer; the parsing actions construct the abstract syntax tree of the input (not the parse tree), which is then pretty-printed.

For tiny and medium-sized grammars the parsers derived in this paper are the implementation of choice: they are fast, at least twice as fast as table-based parsers, and small, the size commensurate to that of table-based ones. For large grammars, the situation is not as clear cut: our parsers are still the fastest; the size of the generated code, however, exceeds that of table-based parsers: the executable is roughly 70% larger. In a sense, the findings mirror the typical time/space trade-off between interpreted and compiled code. It should be noted, however, that the `frown` implementation does not yet incorporate the optimizations described in Section 6.3. The merging of states might ameliorate the size problem.

To summarize, the stackless, table-free parsers are consistently twice as fast as table-based parsers; the larger the grammar the larger the size of the generated parser relative to the table-based one.

As a final remark, it is interesting to note that the use of unsafe features does not necessarily improve the running time. Or to put it

positively, it is pleasing to see that a clean derivation gives rise to an efficient parser.

## 9. Conclusion

We have derived a purely functional implementation of LR parsing. The derived parsers consist of mutually recursive functions corresponding to the states of the underlying LR automaton. Each function takes a number of continuation arguments corresponding to the kernel items of that state. No separate stack is needed as the parsing continuations are immediately applied to the attributes when they become available. Furthermore, the parsers can be typed in a standard Hindley-Milner type system. Polymorphism is only required if one wants to support parsers with multiple entry points.

Once the stage has been set, the derivation proceeds surprisingly smoothly with little or no Eureka steps: the first major transformation (Section 5) is motivated by the need to eliminate the left-recursion present in the initial LL parser; the second major transformation (Section 6.1) eliminates the non-determinism of the intermediate parser, provided the grammar is LR(0).

An interesting feature of our parsers is the passing multiple continuations. These could also be implemented using the multi-function return feature proposed by Shivers (SFar) as an alternative to continuations.

Initial measurements are encouraging: compared to conventional table-based parsers, our parsers are roughly twice as fast. This is not astonishing as our implementation avoids the interpretative overhead of the traditional approach. Furthermore, since our parsers expose the control flow, they are amenable to several optimizations, some of which are specific to the functional representation. On the negative side, the generated parsers are fairly large for realistic grammars. We are confident, however, that the optimizations described in Section 6.3 will ameliorate this problem.

## References

- Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers, Principles, Techniques*. Addison Wesley, 1986.
- Alfred V. Aho and Jeffrey D. Ullman. *The Theory of Parsing, Translation and Compiling*, volume I: Parsing. Prentice-Hall, 1972.
- Lex Augusteijn. *Functional Programming, Program Transformations and Compiler Construction*. PhD thesis, TU Eindhoven, September 1993.



- Olivier Danvy. An extensional characterization of lambda-lifting and lambda-dropping. In Aart Middeldorp and Taisuke Sato, editors, *4th Fuji International Symposium on Functional and Logic Programming (FLOPS'99)*, Tsukuba, Japan, volume 1722, pages 241–250, November 1999.
- Luc Duponcheel and Doaitse Swierstra. A functional program for generating efficient functional LALR(1) parsers, September 2000. unpublished note.
- A. Gill, J. Launchbury, and S.L. Peyton Jones. A short cut to deforestation. In *Proc. Functional Programming Languages and Computer Architecture*, pages 223–232, Copenhagen, 1983. ACM.
- Andy Gill and Simon Marlow. *Happy: The Parser Generator for Haskell, Version 1.15*, 2005. Available from <http://haskell.org/happy/>.
- Ralf Hinze. Constructing tournament representations: An exercise in pointwise relational programming. In Eerke A. Boiten and Bernhard Möller, editors, *Proceedings of the Sixth International Conference on Mathematics of Program Construction (MPC 2002)*, Dagstuhl, Germany, July 8–10, 2002, volume 2386, pages 131–147, July 2002.
- Ralf Hinze. *Frown — an LALR(k) parser generator for Haskell, Version 0.6*, 2005. Available from <http://www.informatik.uni-bonn.de/~ralf/frown/>.
- Donald E. Knuth. On the translation of languages from left to right. *Information and Control*, 8(6):607–639, 1965.
- F.E.J. Kruseman Aretz. On a recursive ascent parser. *Information Processing Letters*, 29:201–206, 1988.
- René Leermakers, Lex Augusteijn, and Frans E.J. Kruseman Aretz. A functional LR-parser. *Theoretical Computer Science*, 104:313–323, 1992.
- René Leermakers. *The Functional Treatment of Parsing*. Kluwer Academic Publishers, July 1993.
- Thomas J. Pennello. Very fast LR parsing. In *SIGPLAN '86: Proceedings of the 1986 SIGPLAN symposium on Compiler construction*, pages 145–151, New York, NY, USA, 1986. ACM Press.
- Peter Pepper. How to obtain powerful parsers that are elegant and practical, March 2004.
- Simon Peyton Jones. *Haskell 98 Language and Libraries*. 2003.
- Francois Pottier and Yann Régis-Gianas. Towards efficient, typed LR parsers. In *ACM SIGPLAN Workshop on ML*, Electronic Notes in Theoretical Computer Science, pages 149–173, September 2005.
- Olin Shivers and David Fisher. Multi-return function call. *Journal of Functional Programming*, to appear.
- Michael Sperber and Peter Thiemann. Generation of LR parsers by partial evaluation. *ACM Transactions on Programming Languages and Systems*, 22(3):224–264, 2000.
- The GHC Team. *The Glorious Glasgow Haskell Compilation System User's Guide, Version 6.4.1*, 2005. Available from <http://www.haskell.org/ghc/>.
- P. Wadler. Deforestation: Transforming programs to eliminate trees. In H. Ganzinger, editor, *Proceedings of the European Symposium on Programming*, volume 300 of *Lecture Notes in Computer Science*, pages 344–358. Springer Verlag, 1988.



$$\begin{aligned}
\text{state}_{T \rightarrow T^*.F} &:: (\text{MulOp} \rightarrow \text{Factor} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{T \rightarrow T^*.F} \text{ k ts} &= \text{case ts of MULOP } v : \text{tr} \rightarrow \text{state}_{T \rightarrow T^*.F} (\text{k } v) \text{ tr} \\
&\quad \quad \quad - \quad \quad \quad \rightarrow \emptyset \\
\text{state}_{T \rightarrow T^*.F} &:: (\text{Factor} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{T \rightarrow T^*.F} \text{ k ts} &= \text{state}_{F \rightarrow \cdot (E)} (\text{reduce}_{F \rightarrow (E)} g_F) \text{ ts} \\
&\quad \cup \text{state}_{F \rightarrow \cdot \text{id}} (\text{reduce}_{F \rightarrow \text{id}} g_F) \text{ ts} \\
\text{where } g_F v &= \text{state}_{T \rightarrow T^*.F} (\text{k } v) \\
\text{state}_{T \rightarrow T^*.F} &:: \mathcal{P} r \rightarrow \mathcal{P} r \\
\text{state}_{T \rightarrow T^*.F} \text{ k ts} &= \text{k ts} \\
\text{state}_{T \rightarrow \cdot F} &:: (\text{Factor} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{T \rightarrow \cdot F} \text{ k ts} &= \text{state}_{F \rightarrow \cdot (E)} (\text{reduce}_{F \rightarrow (E)} g_F) \text{ ts} \\
&\quad \cup \text{state}_{F \rightarrow \cdot \text{id}} (\text{reduce}_{F \rightarrow \text{id}} g_F) \text{ ts} \\
\text{where } g_F v &= \text{state}_{T \rightarrow F} (\text{k } v) \\
\text{state}_{T \rightarrow F} &:: \mathcal{P} r \rightarrow \mathcal{P} r \\
\text{state}_{T \rightarrow F} \text{ k ts} &= \text{k ts} \\
\text{state}_{F \rightarrow \cdot (E)} &:: (\text{Expr} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{F \rightarrow \cdot (E)} \text{ k ts} &= \text{case ts of LPAREN : tr} \rightarrow \text{state}_{F \rightarrow \cdot (E)} \text{ k tr} \\
&\quad \quad \quad - \quad \quad \quad \rightarrow \emptyset \\
\text{state}_{F \rightarrow \cdot (E)} &:: (\text{Expr} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{F \rightarrow \cdot (E)} \text{ k ts} &= \text{state}_{E \rightarrow \cdot E+T} (\text{reduce}_{E \rightarrow E+T} g_E) \text{ ts} \\
&\quad \cup \text{state}_{E \rightarrow \cdot T} (\text{reduce}_{E \rightarrow T} g_E) \text{ ts} \\
\text{where } g_E v &= \text{state}_{F \rightarrow (E \cdot)} (\text{k } v) \\
\text{state}_{F \rightarrow (E \cdot)} &:: \mathcal{P} r \rightarrow \mathcal{P} r \\
\text{state}_{F \rightarrow (E \cdot)} \text{ k ts} &= \text{case ts of RPAREN : tr} \rightarrow \text{state}_{F \rightarrow (E \cdot)} \text{ k tr} \\
&\quad \quad \quad - \quad \quad \quad \rightarrow \emptyset \\
\text{state}_{F \rightarrow (E \cdot)} &:: \mathcal{P} r \rightarrow \mathcal{P} r \\
\text{state}_{F \rightarrow (E \cdot)} \text{ k ts} &= \text{k ts} \\
\text{state}_{F \rightarrow \cdot \text{id}} &:: (\text{String} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{F \rightarrow \cdot \text{id}} \text{ k ts} &= \text{case ts of IDENT } v : \text{tr} \rightarrow \text{state}_{F \rightarrow \cdot \text{id}} (\text{k } v) \text{ tr} \\
&\quad \quad \quad - \quad \quad \quad \rightarrow \emptyset \\
\text{state}_{F \rightarrow \cdot \text{id}} &:: \mathcal{P} r \rightarrow \mathcal{P} r \\
\text{state}_{F \rightarrow \cdot \text{id}} \text{ k ts} &= \text{k ts}
\end{aligned}$$

Note that states with an identical string to the right of the dot are equal:  $\text{state}_{E \rightarrow E+ \cdot T} = \text{state}_{E \rightarrow \cdot T}$  and  $\text{state}_{T \rightarrow T^* \cdot F} = \text{state}_{T \rightarrow \cdot F}$ .

## B. LR(0) Expression Parser

The function  $\text{parse}_E$  is a non-deterministic recursive ascent parser for the expression grammar. The non-determinism is exhibited by the locally defined  $g_X$  functions.

$$\begin{aligned}
\text{parse}_E &:: \mathcal{P} \text{ Expr} \\
\text{parse}_E &= \text{state}_{E^\dagger \rightarrow \cdot E\$} \text{ accept}
\end{aligned}$$

$$\begin{aligned}
\text{state}_{E^\dagger \rightarrow \cdot E\$} &:: (\text{Expr} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{E^\dagger \rightarrow \cdot E\$} \text{ k}_{E^\dagger \rightarrow \cdot E\$} \text{ ts} &= \\
\text{case ts of LPAREN} &: \text{tr} \rightarrow \text{state}_{F \rightarrow (\cdot E)} (\text{reduce}_{F \rightarrow (E)} \text{ g}_F) \text{ tr} \\
\text{IDENT v} &: \text{tr} \rightarrow \text{state}_{F \rightarrow \text{id}} (\text{reduce}_{F \rightarrow \text{id}} \text{ g}_F \text{ v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{where g}_E \text{ v} &= \text{state}_{E^\dagger \rightarrow E \cdot \$} (\text{k}_{E^\dagger \rightarrow \cdot E\$} \text{ v}) \\
&\cup \text{state}_{E \rightarrow E \cdot + T} (\text{reduce}_{E \rightarrow E + T} \text{ g}_E \text{ v}) \\
\text{g}_T \text{ v} &= \text{state}_{E \rightarrow T \cdot} (\text{reduce}_{E \rightarrow T} \text{ g}_E \text{ v}) \\
&\cup \text{state}_{T \rightarrow T \cdot * F} (\text{reduce}_{T \rightarrow T * F} \text{ g}_T \text{ v}) \\
\text{g}_F \text{ v} &= \text{state}_{T \rightarrow F \cdot} (\text{reduce}_{T \rightarrow F} \text{ g}_T \text{ v}) \\
\text{state}_{E^\dagger \rightarrow E \cdot \$} \text{ k ts} &= \\
\text{case ts of EOF} &: \text{tr} \rightarrow \text{k tr} \\
- &\rightarrow \emptyset \\
\text{state}_{E \rightarrow E \cdot + T} \text{ k ts} &= \\
\text{case ts of ADDOP v} &: \text{tr} \rightarrow \text{state}_{E \rightarrow E + \cdot T} (\text{k v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{state}_{E \rightarrow E + \cdot T} &:: (\text{Term} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{E \rightarrow E + \cdot T} \text{ k}_{E \rightarrow E + \cdot T} \text{ ts} &= \\
\text{case ts of LPAREN} &: \text{tr} \rightarrow \text{state}_{F \rightarrow (\cdot E)} (\text{reduce}_{F \rightarrow (E)} \text{ g}_F) \text{ tr} \\
\text{IDENT v} &: \text{tr} \rightarrow \text{state}_{F \rightarrow \text{id}} (\text{reduce}_{F \rightarrow \text{id}} \text{ g}_F \text{ v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{where g}_T \text{ v} &= \text{state}_{E \rightarrow E + T \cdot} (\text{k}_{E \rightarrow E + \cdot T} \text{ v}) \\
&\cup \text{state}_{T \rightarrow T \cdot * F} (\text{reduce}_{T \rightarrow T * F} \text{ g}_T \text{ v}) \\
\text{g}_F \text{ v} &= \text{state}_{T \rightarrow F \cdot} (\text{reduce}_{T \rightarrow F} \text{ g}_T \text{ v}) \\
\text{state}_{E \rightarrow E + T \cdot} \text{ k ts} &= \text{k ts} \\
\text{state}_{E \rightarrow T \cdot} \text{ k ts} &= \text{k ts} \\
\text{state}_{T \rightarrow T \cdot * F} \text{ k ts} &= \\
\text{case ts of MULOP v} &: \text{tr} \rightarrow \text{state}_{T \rightarrow T * \cdot F} (\text{k v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{state}_{T \rightarrow T * \cdot F} &:: (\text{Factor} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{T \rightarrow T * \cdot F} \text{ k}_{T \rightarrow T * \cdot F} \text{ ts} &= \\
\text{case ts of LPAREN} &: \text{tr} \rightarrow \text{state}_{F \rightarrow (\cdot E)} (\text{reduce}_{F \rightarrow (E)} \text{ g}_F) \text{ tr} \\
\text{IDENT v} &: \text{tr} \rightarrow \text{state}_{F \rightarrow \text{id}} (\text{reduce}_{F \rightarrow \text{id}} \text{ g}_F \text{ v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{where g}_F \text{ v} &= \text{state}_{T \rightarrow T * F \cdot} (\text{k}_{T \rightarrow T * \cdot F} \text{ v}) \\
\text{state}_{T \rightarrow T * F \cdot} \text{ k ts} &= \text{k ts} \\
\text{state}_{T \rightarrow F \cdot} \text{ k ts} &= \text{k ts} \\
\text{state}_{F \rightarrow \cdot (E)} \text{ k ts} &= \\
\text{case ts of LPAREN} &: \text{tr} \rightarrow \text{state}_{F \rightarrow (\cdot E)} \text{ k tr} \\
- &\rightarrow \emptyset
\end{aligned}$$

$$\begin{aligned}
\text{state}_{F \rightarrow (\cdot E)} &:: (\text{Expr} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_{F \rightarrow (\cdot E)} k_{F \rightarrow (\cdot E)} \text{ts} &= \\
\quad \text{case ts of LPAREN : tr} &\rightarrow \text{state}_{F \rightarrow (\cdot E)} (\text{reduce}_{F \rightarrow (\cdot E)} g_F) \text{tr} \\
\quad \quad \text{IDENT v : tr} &\rightarrow \text{state}_{F \rightarrow \text{id}} (\text{reduce}_{F \rightarrow \text{id}} g_F v) \text{tr} \\
\quad \quad - &\rightarrow \emptyset \\
\quad \text{where } g_E v &= \text{state}_{E \rightarrow E \cdot + T} (\text{reduce}_{E \rightarrow E \cdot + T} g_E v) \\
&\quad \cup \text{state}_{F \rightarrow (\cdot E)} (k_{F \rightarrow (\cdot E)} v) \\
\quad g_T v &= \text{state}_{E \rightarrow T} (\text{reduce}_{E \rightarrow T} g_E v) \\
&\quad \cup \text{state}_{T \rightarrow T \cdot * F} (\text{reduce}_{T \rightarrow T \cdot * F} g_T v) \\
\quad g_F v &= \text{state}_{T \rightarrow F} (\text{reduce}_{T \rightarrow F} g_T v) \\
\text{state}_{F \rightarrow (\cdot E)} k \text{ts} &= \\
\quad \text{case ts of RPAREN : tr} &\rightarrow \text{state}_{F \rightarrow (E)} k \text{tr} \\
\quad - &\rightarrow \emptyset \\
\text{state}_{F \rightarrow (E)} k \text{ts} &= k \text{ts} \\
\text{state}_{F \rightarrow \cdot \text{id}} k \text{ts} &= \\
\quad \text{case ts of IDENT v : tr} &\rightarrow \text{state}_{F \rightarrow \text{id}} (k v) \text{tr} \\
\quad - &\rightarrow \emptyset \\
\text{state}_{F \rightarrow \text{id}} k \text{ts} &= k \text{ts}
\end{aligned}$$

### C. LALR(1) Expression Parser

The function  $\text{parse}_E$  is a deterministic recursive ascent parser for the expression grammar corresponding to a conventional LALR(1) parser.

$$\begin{aligned}
\text{parse}_E &:: \mathcal{P} \text{Expr} \\
\text{parse}_E &= \text{state}_1 \text{ accept} \\
\text{state}_1 &:: (\text{Expr} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_1 k_{E \rightarrow \cdot E\$} \text{ts} &= \\
\quad \text{case ts of LPAREN : tr} &\rightarrow \text{state}_{11} (\text{reduce}_{F \rightarrow (\cdot E)} g_F) \text{tr} \\
\quad \quad \text{IDENT v : tr} &\rightarrow \text{state}_{13} (\text{reduce}_{F \rightarrow \text{id}} g_F v) \text{tr} \\
\quad \quad - &\rightarrow \emptyset \\
\quad \text{where } g_E v &= \text{state}_2 (k_{E \rightarrow \cdot E\$} v) (\text{reduce}_{E \rightarrow E \cdot + T} g_E v) \\
\quad g_T v &= \text{state}_7 (\text{reduce}_{E \rightarrow T} g_E v) (\text{reduce}_{T \rightarrow T \cdot * F} g_T v) \\
\quad g_F v &= \text{state}_{10} (\text{reduce}_{T \rightarrow F} g_T v) \\
\text{state}_2 &:: \mathcal{P} r \rightarrow (\text{AddOp} \rightarrow \text{Term} \rightarrow \mathcal{P} r) \rightarrow \mathcal{P} r \\
\text{state}_2 k_{E \rightarrow \cdot E\$} k_{E \rightarrow E \cdot + T} \text{ts} &= \\
\quad \text{case ts of ADDOP v : tr} &\rightarrow \text{state}_5 (k_{E \rightarrow E \cdot + T} v) \text{tr} \\
\quad \quad \text{EOF : tr} &\rightarrow \text{state}_3 k_{E \rightarrow \cdot E\$} \text{tr} \\
\quad \quad - &\rightarrow \emptyset \\
\text{state}_3 &:: \mathcal{P} r \rightarrow \mathcal{P} r \\
\text{state}_3 k_{E \rightarrow \cdot E\$} \text{ts} &= k_{E \rightarrow \cdot E\$} \text{ts}
\end{aligned}$$

$$\begin{aligned}
\text{state}_4 &:: (\text{AddOp} \rightarrow \text{Term} \rightarrow \mathcal{P} \text{ r}) \rightarrow \mathcal{P} \text{ r} \rightarrow \mathcal{P} \text{ r} \\
\text{state}_4 \text{ k}_{E \rightarrow E \cdot + T} \text{ k}_{F \rightarrow (E \cdot)} \text{ ts} &= \\
\text{case ts of RPAREN : tr} &\rightarrow \text{state}_{12} \text{ k}_{F \rightarrow (E \cdot)} \text{ tr} \\
\text{ADDOP v : tr} &\rightarrow \text{state}_5 (\text{k}_{E \rightarrow E \cdot + T} \text{ v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{state}_5 &:: (\text{Term} \rightarrow \mathcal{P} \text{ r}) \rightarrow \mathcal{P} \text{ r} \\
\text{state}_5 \text{ k}_{E \rightarrow E \cdot + T} \text{ ts} &= \\
\text{case ts of LPAREN : tr} &\rightarrow \text{state}_{11} (\text{reduce}_{F \rightarrow (E)} \text{ g}_F) \text{ tr} \\
\text{IDENT v : tr} &\rightarrow \text{state}_{13} (\text{reduce}_{F \rightarrow \text{id}} \text{ g}_F \text{ v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{where g}_T \text{ v} &= \text{state}_6 (\text{k}_{E \rightarrow E \cdot + T} \text{ v}) (\text{reduce}_{T \rightarrow T * F} \text{ g}_T \text{ v}) \\
\text{g}_F \text{ v} &= \text{state}_{10} (\text{reduce}_{T \rightarrow F} \text{ g}_T \text{ v}) \\
\text{state}_6 &:: \mathcal{P} \text{ r} \rightarrow (\text{MulOp} \rightarrow \text{Factor} \rightarrow \mathcal{P} \text{ r}) \rightarrow \mathcal{P} \text{ r} \\
\text{state}_6 \text{ k}_{E \rightarrow E \cdot + T} \text{ k}_{T \rightarrow T \cdot * F} \text{ ts} &= \\
\text{case ts of MULOP v : tr} &\rightarrow \text{state}_8 (\text{k}_{T \rightarrow T \cdot * F} \text{ v}) \text{ tr} \\
- &\rightarrow \text{k}_{E \rightarrow E \cdot + T} \text{ ts} \\
\text{state}_7 &:: \mathcal{P} \text{ r} \rightarrow (\text{MulOp} \rightarrow \text{Factor} \rightarrow \mathcal{P} \text{ r}) \rightarrow \mathcal{P} \text{ r} \\
\text{state}_7 \text{ k}_{E \rightarrow T} \text{ k}_{T \rightarrow T \cdot * F} \text{ ts} &= \\
\text{case ts of MULOP v : tr} &\rightarrow \text{state}_8 (\text{k}_{T \rightarrow T \cdot * F} \text{ v}) \text{ tr} \\
- &\rightarrow \text{k}_{E \rightarrow T} \text{ ts} \\
\text{state}_8 &:: (\text{Factor} \rightarrow \mathcal{P} \text{ r}) \rightarrow \mathcal{P} \text{ r} \\
\text{state}_8 \text{ k}_{T \rightarrow T \cdot * F} \text{ ts} &= \\
\text{case ts of LPAREN : tr} &\rightarrow \text{state}_{11} (\text{reduce}_{F \rightarrow (E)} \text{ g}_F) \text{ tr} \\
\text{IDENT v : tr} &\rightarrow \text{state}_{13} (\text{reduce}_{F \rightarrow \text{id}} \text{ g}_F \text{ v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{where g}_F \text{ v} &= \text{state}_9 (\text{k}_{T \rightarrow T \cdot * F} \text{ v}) \\
\text{state}_9 &:: \mathcal{P} \text{ r} \rightarrow \mathcal{P} \text{ r} \\
\text{state}_9 \text{ k}_{T \rightarrow T \cdot * F} \text{ ts} &= \text{k}_{T \rightarrow T \cdot * F} \text{ ts} \\
\text{state}_{10} &:: \mathcal{P} \text{ r} \rightarrow \mathcal{P} \text{ r} \\
\text{state}_{10} \text{ k}_{T \rightarrow F} \text{ ts} &= \text{k}_{T \rightarrow F} \text{ ts} \\
\text{state}_{11} &:: (\text{Expr} \rightarrow \mathcal{P} \text{ r}) \rightarrow \mathcal{P} \text{ r} \\
\text{state}_{11} \text{ k}_{F \rightarrow (E)} \text{ ts} &= \\
\text{case ts of LPAREN : tr} &\rightarrow \text{state}_{11} (\text{reduce}_{F \rightarrow (E)} \text{ g}_F) \text{ tr} \\
\text{IDENT v : tr} &\rightarrow \text{state}_{13} (\text{reduce}_{F \rightarrow \text{id}} \text{ g}_F \text{ v}) \text{ tr} \\
- &\rightarrow \emptyset \\
\text{where g}_E \text{ v} &= \text{state}_4 (\text{reduce}_{E \rightarrow E \cdot + T} \text{ g}_E \text{ v}) (\text{k}_{F \rightarrow (E)} \text{ v}) \\
\text{g}_T \text{ v} &= \text{state}_7 (\text{reduce}_{E \rightarrow T} \text{ g}_E \text{ v}) (\text{reduce}_{T \rightarrow T * F} \text{ g}_T \text{ v}) \\
\text{g}_F \text{ v} &= \text{state}_{10} (\text{reduce}_{T \rightarrow F} \text{ g}_T \text{ v}) \\
\text{state}_{12} &:: \mathcal{P} \text{ r} \rightarrow \mathcal{P} \text{ r} \\
\text{state}_{12} \text{ k}_{F \rightarrow (E)} \text{ ts} &= \text{k}_{F \rightarrow (E)} \text{ ts} \\
\text{state}_{13} &:: \mathcal{P} \text{ r} \rightarrow \mathcal{P} \text{ r} \\
\text{state}_{13} \text{ k}_{F \rightarrow \text{id}} \text{ ts} &= \text{k}_{F \rightarrow \text{id}} \text{ ts}
\end{aligned}$$

Note that  $\text{state}_6$  and  $\text{state}_7$  are identical since the two sets of kernel items do not differ in the strings to the right of the dot.

#### D. Smaller LALR(1) Expression Parser

Again,  $\text{parse}_E$  is a deterministic recursive ascent parser for the expression grammar. Compared to the parser in Appendix C, trivial states of the form  $\text{state}_i$   $k_q = k_q$  have been inlined. Furthermore, common actions have been factored out:  $\text{state}_1$  is defined in terms of  $\text{state}_5$  which in turn depends on  $\text{state}_8$ .

```

 $\text{parse}_E$                 ::  $\mathcal{P}$  Expr
 $\text{parse}_E$                 =  $\text{state}_1$  accept
 $\text{state}_1$                 :: (Expr  $\rightarrow$   $\mathcal{P}$  r)  $\rightarrow$   $\mathcal{P}$  r
 $\text{state}_1$   $k_{E \rightarrow \cdot E \$}$ 
  where  $g_E$  v          =  $\text{state}_5$  ( $\text{reduce}_{E \rightarrow T}$   $g_E$ )
                    =  $\text{state}_2$  ( $k_{E \rightarrow \cdot E \$}$  v) ( $\text{reduce}_{E \rightarrow E+T}$   $g_E$  v)
 $\text{state}_2$                 ::  $\mathcal{P}$  r  $\rightarrow$  (AddOp  $\rightarrow$  Term  $\rightarrow$   $\mathcal{P}$  r)  $\rightarrow$   $\mathcal{P}$  r
 $\text{state}_2$   $k_{E \rightarrow E \cdot \$}$   $k_{E \rightarrow E \cdot +T}$  ts =
  case ts of ADDOP v : tr  $\rightarrow$   $\text{state}_5$  ( $k_{E \rightarrow E \cdot +T}$  v) tr
             EOF      : tr  $\rightarrow$   $k_{E \rightarrow E \cdot \$}$  tr
             -        : tr  $\rightarrow$   $\emptyset$ 
 $\text{state}_4$                 :: (AddOp  $\rightarrow$  Term  $\rightarrow$   $\mathcal{P}$  r)  $\rightarrow$   $\mathcal{P}$  r  $\rightarrow$   $\mathcal{P}$  r
 $\text{state}_4$   $k_{E \rightarrow E \cdot +T}$   $k_{F \rightarrow (E \cdot)}$  ts =
  case ts of RPAREN : tr  $\rightarrow$   $k_{F \rightarrow (E \cdot)}$  tr
             ADDOP v : tr  $\rightarrow$   $\text{state}_5$  ( $k_{E \rightarrow E \cdot +T}$  v) tr
             -        : tr  $\rightarrow$   $\emptyset$ 
 $\text{state}_5$                 :: (Term  $\rightarrow$   $\mathcal{P}$  r)  $\rightarrow$   $\mathcal{P}$  r
 $\text{state}_5$   $k_{E \rightarrow E+ \cdot T}$ 
  where  $g_T$  v          =  $\text{state}_8$  ( $\text{reduce}_{T \rightarrow F}$   $g_T$ )
                    =  $\text{state}_6$  ( $k_{E \rightarrow E+ \cdot T}$  v) ( $\text{reduce}_{T \rightarrow T * F}$   $g_T$  v)
 $\text{state}_6$                 ::  $\mathcal{P}$  r  $\rightarrow$  (MulOp  $\rightarrow$  Factor  $\rightarrow$   $\mathcal{P}$  r)  $\rightarrow$   $\mathcal{P}$  r
 $\text{state}_6$   $k_{E \rightarrow E+T \cdot}$   $k_{T \rightarrow T \cdot * F}$  ts =
  case ts of MULOP v : tr  $\rightarrow$   $\text{state}_8$  ( $k_{T \rightarrow T \cdot * F}$  v) tr
             -        : tr  $\rightarrow$   $k_{E \rightarrow E+T \cdot}$  ts
 $\text{state}_8$                 :: (Factor  $\rightarrow$   $\mathcal{P}$  r)  $\rightarrow$   $\mathcal{P}$  r
 $\text{state}_8$   $k_{T \rightarrow T * \cdot F}$  ts =
  case ts of LPAREN : tr  $\rightarrow$   $\text{state}_{11}$  ( $\text{reduce}_{F \rightarrow (E)}$   $k_{T \rightarrow T * \cdot F}$ ) tr
             IDENT v : tr  $\rightarrow$   $\text{reduce}_{F \rightarrow \text{id}}$   $k_{T \rightarrow T * \cdot F}$  v tr
             -        : tr  $\rightarrow$   $\emptyset$ 
 $\text{state}_{11}$             :: (Expr  $\rightarrow$   $\mathcal{P}$  r)  $\rightarrow$   $\mathcal{P}$  r
 $\text{state}_{11}$   $k_{F \rightarrow (\cdot E)}$ 
  where  $g_E$  v          =  $\text{state}_5$  ( $\text{reduce}_{E \rightarrow T}$   $g_E$ )
                    =  $\text{state}_4$  ( $\text{reduce}_{E \rightarrow E+T}$   $g_E$  v) ( $k_{F \rightarrow (\cdot E)}$  v)

```

