

Function `lfp` applies a function `f` to an argument `x` until it reaches a fixed point. In order to determine whether an argument is a fixed point, `f x` is compared with `x`. In the case of grammar analysis problems, the first components of the elements of the grammar analysis problems are always the nonterminals of the given grammar. It follows that the first components of the elements of `f x` and `x` are always equal, and that equality depends just on the second components of the problems. The first condition of function `lfp` may be replaced by `map snd x = map snd (f x)`.

## 6 Conclusions

This paper discusses bottom-up grammar analysis problems. We give a very general specification of bottom-up grammar analysis problems, and from this specification we derive, by means of program transformation applying laws to the components of the intermediate expressions, an algorithm for performing bottom-up grammar analysis. The driving force in the derivation of the algorithm is the construction of a fixed point. To obtain such a fixed point a number of conditions have to be imposed upon the components of the bottom-up grammar analysis problem. Thus we derive both the algorithm and the conditions under which the fixed point exists in one go. The derivation is an example of a derivation of a real-world program, which would have been difficult to obtain without a derivation. The research reported on in this paper is still in progress: in the next version we want to split the calculation in two parts. The first part of the derivation assumes that the function that computes the property of a parse tree is a *Rosetree* catamorphism and the second part of the derivation adds, if necessary, the extra information (for example in the case of *firsts*, where we use information about the *empties*). This simplifies the derivation. Future research will be directed towards the derivation of an algorithm for top-down grammar analysis.

## References

- [1] R.S. Bird. An introduction to the theory of lists. In M. Broy, editor, *Logic of Programming and Calculi of Discrete Design*, volume F36 of *NATO ASI Series*, pages 5–42. Springer-Verlag, 1987.
- [2] B.A. Davey and H.A. Priestley. *Introduction to Lattices and Order*. Cambridge University Press, 1990.
- [3] J.H. Fasel, P. Hudak, S. Peyton Jones, and P. Wadler. Sigplan Notices Special Issue on the Functional Programming Language Haskell. *ACM SIGPLAN notices*, 27(5), 1992.
- [4] J. Jeuring. *Theories for Algorithm Calculation*. PhD thesis, Utrecht University, 1993.
- [5] M.P. Jones. Introduction to Gofer 2.20. Programming Research Group, Oxford University, 1992.
- [6] G. Malcolm. Data structures and program transformation. *Science of Computer Programming*, 14:255–279, 1990.
- [7] L. Meertens. Algorithmics—towards programming as a mathematical activity. In J.W. de Bakker, M. Hazewinkel, and J.K. Lenstra, editors, *Proceedings of the CWI Symposium on Mathematics and Computer Science*, volume 1 of *CWI Monographs*, pages 289–334. North-Holland, 1986.
- [8] Torben Mogensen. Ratatosk – a parser generator and scanner generator for Gofer. Published on `comp.lang.functional`, 1993.
- [9] Ulrich Möncke and Reinhard Wilhelm. Grammar flow analysis. In *Attribute Grammars, Applications and Systems, SAGA '91*, pages 151–186. Springer-Verlag, New York, 1991. LNCS 545.
- [10] Simon L. Peyton Jones. Yacc in Sasl – an exercise in functional programming. *Software-Practice and Experience*, 15(8):807–820, 1985.
- [11] P. Wadler. Comprehending monads. *Mathematical Structures in Computer Science*, 2:461–493, 1992.

$$\begin{aligned}
& (([] =) \cdot \# / \cdot \text{sen}^*) x \\
= & ( [] =) \cdot \# / = \text{and} \cdot ( [] =)^*, (7) \\
& (\text{and} \cdot (([] =) \cdot \text{sen})^*) x \\
= & \text{definition of } qn \text{ (see below), and } p \\
& qn \text{ nt } (p^* x)
\end{aligned}$$

Function  $qn$  is defined by

$$qn \text{ nt } x = \text{and } x$$

Furthermore, we have to show that  $\wedge$ , the operator of the reduction for  $\text{and}$ , distributes over  $\vee$ , that is,

$$\begin{aligned}
a \wedge (b \vee c) &= (a \wedge b) \vee (a \wedge c) \\
(a \vee b) \wedge c &= (a \wedge c) \vee (b \wedge c)
\end{aligned}$$

and that  $\text{false}$  is a zero of  $\wedge$ . These equalities hold for  $\text{false}$ ,  $\vee$  and  $\wedge$ . Finally, we have to show that  $\wedge$  is monotonic in both arguments. This requirement is satisfied too.

FIRST

We verify the conditions the components of the definition of the bottom-up grammar analysis problem FIRST have to satisfy.

First the join semilattice  $(c, \sqsubseteq_c)$  upon which the join semilattice  $(E, \sqsubseteq_E)$  is built is the join semilattice of terminals, where  $c$  is the set of terminals, the relation  $\sqsubseteq_c$  is the subset relation,  $[]$  is the bottom of  $c$ , and the join  $\sqcup_c$  is set union, or  $nub \cdot \#$ . Clearly, set union is associative, and  $[]$  is the unit of set union.

For the second assumption, we have to construct a function  $H$  that can be implemented as an efficient program, such that equation (22) holds. The condition (26) given in the previous subsection does not hold for function  $p$  defined by

$$p = \text{take } 1 \cdot \text{sen}$$

It is not difficult to find a *Rosetree* catamorphism for  $p$ , so (25) is satisfied, but the second requirement (26) does not hold. It follows that we have to find another way to construct function  $H$ . Function  $H$  is defined by

$$H \text{ nt} = \text{foldr } t \ 1_{\oplus}$$

where function  $t$  is defined as follows. If the current symbol in the right-hand side of a production is a terminal, then the symbols that can appear as the first symbol of a string are the symbols found until then, and no more.

$$t (T \ b, y) \ x = y$$

If the current symbol in the right-hand side of a production is a nonterminal  $N \ a$ , then we distinguish two cases depending on whether or not  $N \ a$  can derive the empty string. If  $N \ a$  can derive the empty string, then the symbols that can appear as the first symbol of a string are the symbols found until then together with the first symbols of the remaining part of the production. If  $N \ a$  cannot derive the empty string then the symbols that can appear as the first symbol of a string are the symbols found until then, and no more.

$$\begin{aligned}
& t (N \ a, y) \ x \\
= & \begin{cases} nub (y \# x) & \text{if } \text{at } \text{empties } a \\ y & \text{otherwise} \end{cases}
\end{aligned}$$

We can prove equation (22) for function  $H$  thus defined by induction to the structure of lists: apply both sides to  $[]$  and  $[a] \# x$ , and show that the resulting expressions have the same recursive structure. The proof of the fact that  $H$  is monotonic is easy and omitted.

## 5.6 Implementation

The definitions given in the previous subsections are translated into Gofer as follows. Some rather obvious alterations of these functions increase the efficiency of the program. These alterations are discussed after the following program.

```

ag g pn pl
= lfp
  k
  (map (split id (\x -> bottom)) nt's)
  where
  nt's = nts g
  k x = map (split id (j x)) nt's
  j x nt = (r x nt)
          'join'
          ((lub
            .map (pn nt . map (r' x))
            .rhss g)
           nt
          )
  r x nt = at x nt
  r' x (N a) = (N a, r x a)
  r' x (T b) = (T b, pl b)

```

We discuss one of the possibly many ways in which a more efficient program can be obtained.

mation, and function  $H$  is therefore highly inefficient. However, it follows from this definition that the only condition that has to be satisfied in order to solve a bottom-up grammar analysis problem is the first condition given above. To obtain a practical solution for a bottom-up grammar analysis problem we discuss a special case in which we can find a monotonic function  $H$  that can be implemented as an efficient program.

Suppose the property function  $p$  is a catamorphism on *Rosetree*. Then we have for function  $pn\ nt$  satisfying assumption (21)

$$pn\ nt = qn\ nt \cdot exr* \quad (24)$$

where function  $qn$  is the function of the *Rosetree* catamorphism for  $p$ . For the left-hand expression of equation (22) we now calculate as follows.

$$\begin{aligned} & \oplus/ \cdot (pn\ nt)* \cdot cp \cdot zri* \\ = & \quad (24) \\ & \oplus/ \cdot (qn\ nt \cdot exr*)* \cdot cp \cdot zri* \\ = & \quad \text{map distributivity (7)} \\ & \oplus/ \cdot (qn\ nt)* \cdot exr** \cdot cp \cdot zri* \\ = & \quad (14), \text{ map distributivity (7)} \\ & \oplus/ \cdot (qn\ nt)* \cdot cp \cdot (exr* \cdot zri)* \\ = & \quad \text{product/zip calculation (omitted)} \\ & \oplus/ \cdot (qn\ nt)* \cdot cp \cdot exr* \\ = & \quad \text{assume equation (26) below} \\ & qn\ nt \cdot (\oplus/)* \cdot exr* \\ = & \quad \text{map distributivity (7), (3), (7)} \\ & qn\ nt \cdot exr* \cdot (id \times \oplus/)* \end{aligned}$$

It follows that if we assume that there exists a function  $qn$  such that

$$p \cdot (Node\ nt) = qn\ nt \cdot p* \quad (25)$$

$$\oplus/ \cdot (qn\ nt)* \cdot cp = qn\ nt \cdot (\oplus/)* \quad (26)$$

then function  $H$  can be defined by

$$H\ nt = qn\ nt \cdot exr*$$

The second assumption is still rather unwieldy, and can be simplified. To obtain a simpler condition we apply the theory for  $cp$  developed in [4]. For that purpose, we first assume that function  $qn\ nt$  is a reduction, that is, there exists an operator  $\otimes$  with unit  $u$  such that

$$qn\ nt = \otimes/$$

Now we apply a theorem from [4], which states that (26) holds, provided the sections  $(a \otimes)$  and  $(\otimes a)$  distribute over operator  $\oplus$ .

$$\begin{aligned} a \otimes (b \oplus c) &= (a \otimes b) \oplus (a \otimes c) \\ (a \oplus b) \otimes c &= (a \otimes c) \oplus (b \otimes c) \end{aligned}$$

and provided for all  $y$ ,  $1_{\oplus} \otimes y = y \otimes 1_{\oplus} = 1_{\oplus}$ . Function  $H\ nt$  is monotonic provided function  $qn\ nt$  is monotonic, and function  $qn\ nt$  is monotonic provided operator  $\otimes$  is monotonic in both arguments.

## 5.5 Examples

This section shows how we apply the theory derived in the previous section to the examples of bottom-up grammar analysis problems given in Section 4. The algorithm derived in the previous section can be used to solve a bottom-up grammar analysis problem provided the components of the grammar analysis problem satisfy the conditions given in the previous section.

### EMPTY

We verify the conditions the components of the definition of the bottom-up grammar analysis problem EMPTY have to satisfy.

First, the join semilattice  $(c, \sqsubseteq_c)$  upon which the join semilattice  $(E, \sqsubseteq_E)$  is built is the join semilattice of booleans, where  $c$  is the set  $\{true, false\}$ , the relation  $\sqsubseteq_c$  is defined by  $false \sqsubseteq_c true$ ,  $false$  is the bottom of  $c$ , and the join  $\sqcup_c$  is the operator  $\vee$ . Clearly,  $\vee$  is associative, and  $false$  is the unit of  $\vee$ .

For the second assumption, we have to construct a function  $H$  such that equation (22) holds. To obtain a definition of function  $H$  that can be implemented as an efficient program, we verify the conditions listed in the previous subsection. We have to show that function  $p$  defined by

$$p = ([ ] =) \cdot sen$$

is a *Rosetree* catamorphism, i.e., there should exist a function  $qn$  such that

$$p (Node\ nt\ x) = qn\ nt (p* x)$$

A definition of function  $qn$  is obtained as follows.

$$\begin{aligned} & p (Node\ nt\ x) \\ = & \quad \text{definition of } h \\ & (([ ] =) \cdot sen \cdot Node\ nt) x \\ = & \quad \text{definition of } sen \end{aligned}$$

function  $J$  satisfying equation (18). Function  $J$  is defined by

$$\begin{aligned} & J \ x \ a \\ = & \\ & (r \ x \ a) \oplus \\ & ((\oplus / \cdot (H \ a \cdot (r' \ x)*)*) (rhss \ g \ a)) \end{aligned} \quad (23)$$

It remains to prove that there exists a function  $H$  such that equation (22) is satisfied, and that function  $K$  defined in equation (19) is monotonic. The latter condition is discussed in the following subsection, and the former condition in the subsection thereafter.

We give an operational interpretation of the functions we have derived. Given a grammar  $g$  and a CPO  $(E, \sqsubseteq_E)$ , we compute the least fixed point of function  $K$ , starting with  $K \perp$ , where  $\perp$  is the bottom of  $E$ , and repeatedly applying  $K$  until we find a value  $x$  such that  $K \ x = x$ . Function  $K$  applies function  $J$  to all nonterminals of  $g$ . Function  $J$  takes the old value of  $K$  and a nonterminal  $nt$ , and returns the new value for  $nt$  by applying the function  $H \ nt \cdot (r' \ x)^*$  to all right-hand sides of the productions of nonterminal  $nt$ . The results are combined by taking the join  $\oplus /$  of the values thus obtained, and, finally, by joining the result with the old value for  $nt$ .

### 5.3 Function $K$ is monotonic

In order to guarantee the existence of the least fixed point of function  $K : E \rightarrow E$  defined by

$$K \ x = (id \ \Delta \ (J \ x))^* \ (nts \ g)$$

where function  $J$  is defined in equation (23), we have to show that function  $K$  is monotonic, i.e., for  $a, a' \in E$ ,  $K$  has to satisfy

$$a \sqsubseteq_E a' \Rightarrow K \ a \sqsubseteq_E K \ a'$$

It is easily verified that function  $K$  is monotonic, provided function  $J$  is monotonic on  $E$  in its first argument, that is, provided

$$a \sqsubseteq_E a' \Rightarrow J \ a \ nt \sqsubseteq_E J \ a' \ nt$$

Function  $J$  is monotonic in  $a$ , provided operator  $\oplus$  is monotonic in both its arguments, and function  $H$  satisfying equation (22) is monotonic in its second argument. These are the last conditions we impose upon the components of a bottom-up grammar analysis problem. An example of an operator  $\oplus$  that is monotonic in both its arguments is

the operator join  $\sqcup_c$  of the semilattice  $c$  by means of which the semilattice  $E$  is defined. In the examples of the following subsection and the program in subsection 5.6, operator  $\oplus$  will be the join of a join semilattice, which renders the verification of monotonicity of  $\oplus$  trivial. Since  $K$  is monotonic and the domain of  $K$  is finite, the least fixed point of  $K$  can be found in finite time.

### 5.4 The conditions

In the previous subsections we have derived a function  $K$  by means of which a bottom-up grammar analysis problem can be solved. In the derivation we have imposed a number of conditions upon the components of the grammar analysis problems. This subsection discusses these conditions.

The first condition we imposed upon bottom-up grammar analysis problems is the following. We suppose there exists a join semilattice  $(c, \sqsubseteq_c)$  such that  $1_\oplus$  is the bottom of  $c$ , and  $\oplus$  is the join  $\sqcup_c$  of  $c$ .

For the second condition we suppose that there exists a monotonic function  $H$ , such that the following equality holds.

$$\begin{aligned} & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot zri^* \\ = & \\ & H \ nt \cdot (id \ \times \ \oplus /)^* \end{aligned}$$

There exists a trivial but rather useless monotonic function  $H$  such that the above equality is satisfied:

$$\begin{aligned} & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot zri^* \\ = & \quad (1) \\ & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot zri^* \cdot (exl \ \cdot id \ \Delta \ \oplus /)^* \\ = & \quad \text{map-distributivity (7)} \\ & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot zri^* \cdot exl^* \cdot (id \ \Delta \ \oplus /)^* \\ = & \quad \text{definition of } H \text{ below} \\ & H \ nt \cdot (id \ \Delta \ \oplus /)^* \end{aligned}$$

where function  $H$  is defined by

$$\begin{aligned} & H \ nt \\ = & \\ & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot zri^* \cdot exl^* \end{aligned}$$

Function  $H$  recomputes the required information from scratch instead of using the available infor-

and it follows that if we define function  $r$  by

$$r \ x \ a \ = \ at \ x \ a$$

where function  $at$  is defined by

$$at \ x \ a \ = \ head \ [y \ | \ (a, y) \leftarrow x]$$

then

$$af \cdot gh \ n \cdot N \ = \ r \ (agn \ n) \quad (20)$$

This equation is used to express the left-hand argument of operator  $\oplus$  in terms of  $agn \ n$ . It remains to express the right-hand argument of operator  $\oplus$  in terms of  $agn \ n$ . We calculate as follows for  $af \ rh$

$$\begin{aligned} & af \ rh \\ = & \text{definition of } rh \\ & af \ [Node \ nt \ c \ | \ rhs \leftarrow rhss \ g \ nt, \\ & \quad \quad \quad c \leftarrow cp \ ((gh \ n)^* \ rhs)] \\ = & \text{calculation from Section 4} \\ & (af \cdot (Node \ nt)^* \cdot \# / \cdot cp^* \cdot (gh \ n)^{**}) \\ & (rhss \ g \ nt) \end{aligned}$$

If we can push  $af$  to the right within the map  $(gh \ n)^{**}$  in the composition of functions of the last expression in the above calculation, then we can use equation (20) again to obtain an expression of the desired form. Aiming at pushing  $af$  to the right then, we proceed with the composition of functions  $af \cdot (Node \ nt)^* \cdot \# / \cdot cp^* \cdot (gh \ n)^{**}$ . Abbreviate function  $Node \ nt$  to  $mt$ .

$$\begin{aligned} & af \cdot mt^* \cdot \# / \cdot cp^* \cdot (gh \ n)^{**} \\ = & \text{definition of } af \\ & \oplus / \cdot p^* \cdot mt^* \cdot \# / \cdot cp^* \cdot (gh \ n)^{**} \\ = & \text{map-distributivity (7)} \\ & \oplus / \cdot (p \cdot mt)^* \cdot \# / \cdot cp^* \cdot (gh \ n)^{**} \\ = & \text{equation (12)} \\ & \oplus / \cdot \# / \cdot (p \cdot mt)^{**} \cdot cp^* \cdot (gh \ n)^{**} \\ = & \text{equation (13)} \\ & \oplus / \cdot \oplus / * \cdot (p \cdot mt)^{**} \cdot cp^* \cdot (gh \ n)^{**} \\ = & \text{map-distributivity (7)} \\ & \oplus / \cdot (\oplus / \cdot (p \cdot mt)^* \cdot cp \cdot (gh \ n)^*)^* \end{aligned}$$

At this point of the calculation we assume that there exists a function  $pn$  such that

$$p \ (Node \ nt \ x) \ = \ pn \ nt \ ((top \ \Delta \ p)^* \ x) \quad (21)$$

This condition is not unreasonable: for all *Rose-Tree* catamorphisms there exists such a function  $pn$ . We proceed the calculation with the expression within the map in the last expression of the above calculation.

$$\begin{aligned} & \oplus / \cdot (p \cdot mt)^* \cdot cp \cdot (gh \ n)^* \\ = & \text{assumption (21)} \\ & \oplus / \cdot (pn \ nt \cdot (top \ \Delta \ p)^*)^* \cdot cp \cdot (gh \ n)^* \\ = & \text{map-distributivity (7)} \\ & \oplus / \cdot (pn \ nt)^* \cdot (top \ \Delta \ p)^{**} \cdot cp \cdot (gh \ n)^* \\ = & (14) \\ & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot (top \ \Delta \ p)^{**} \cdot (gh \ n)^* \\ = & \text{map-distributivity (7)} \\ & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot ((top \ \Delta \ p)^* \cdot (gh \ n))^* \\ = & \text{introduction of function } zri \ \text{below} \\ & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot zri^* \cdot (id \ \Delta \ (p^* \cdot gh \ n))^* \\ = & \text{assume equation (22) below} \\ & H \ nt \cdot (id \ \times \ \oplus /)^* \cdot (id \ \Delta \ (p^* \cdot gh \ n))^* \\ = & \text{map-distributivity (7), (5)} \\ & H \ nt \cdot (id \ \Delta \ (\oplus / \cdot p^* \cdot gh \ n))^* \\ = & \text{introduction of } r' \ \text{below} \\ & H \ nt \cdot (r' \ (agn \ n))^* \end{aligned}$$

In this calculation we have assumed the existence of three functions:  $zri$ ,  $H$ , and  $r'$  such that a number of properties is satisfied. Function  $zri$  is defined by

$$zri \ = \ \Upsilon \cdot repeat \ \times \ id$$

where function  $repeat$  takes an element  $a$ , and returns an infinite list of  $a$ 's. We omit the proof of the fact that function  $zri$  satisfies the following equality.

$$zri \cdot id \ \Delta \ (p^* \cdot gh \ n) \ = \ (top \ \Delta \ p)^* \cdot (gh \ n)$$

Furthermore, we have assumed the existence of a function  $H$  such that the following equality is satisfied.

$$\begin{aligned} & \oplus / \cdot (pn \ nt)^* \cdot cp \cdot zri^* \\ = & \hspace{15em} (22) \\ & H \ nt \cdot (id \ \times \ \oplus /)^* \end{aligned}$$

Finally, function  $r'$  is defined by

$$\begin{aligned} r' \ x \ (N \ a) & \ = \ (N \ a, r \ x \ a) \\ r' \ x \ (T \ b) & \ = \ (T \ b, pl \ b) \end{aligned}$$

The above derivation shows that if there exist a function  $H$  satisfying (22), then there exists a

are all equal to  $1_{\oplus}$ . This suggests to construct the following join semilattice. Let  $E$  be the set of lists  $x$  of length equal to the number of nonterminals of  $g$  of which  $exl * x = nts g$ , and of which the second component of each element is an element of  $c$ , the result type of operator  $\oplus$ .

For the definition of the relation  $\sqsubseteq_E$  and the join  $\sqcup_E$ , we suppose that there exists a relation  $\sqsubseteq_c$  such that  $(c, \sqsubseteq_c)$  with join  $\sqcup_c$  is a join semilattice, and such that the unit  $1_{\oplus}$  of operator  $\oplus$  occurring in the definition of a bottom-up grammar analysis problem is the bottom of  $c$ .

Both the relation  $\sqsubseteq_E$  and the join  $\sqcup_E$  are now straightforward extensions of  $\sqsubseteq_c$  and  $\sqcup_c$ , respectively. Relation  $\sqsubseteq_E$  is defined by pairwise comparing elements with  $\sqsubseteq_c$ .

$$\begin{aligned} & x \sqsubseteq_E y \\ \equiv & \\ & \text{and } (exr * x \curlywedge_{\sqsubseteq_c} exr * y) \end{aligned}$$

where function *and* is the reduce  $\wedge$ , and where  $\curlywedge_{\oplus}$ , with  $\oplus$  a binary function, zips two lists of equal length to a list of pairs, and then applies operator  $\oplus$  to all pairs in the list. The join of two elements is defined by pairwise joining the second components of the pairs.

$$\begin{aligned} & x \sqcup_E y \\ = & \\ & exl * x \curlywedge (exr * x \curlywedge_{\sqcup_c} exr * y) \end{aligned}$$

It is easy to prove that  $agn\ 0$  is the bottom of  $E$ , using the fact that  $1_{\oplus}$  is the bottom of  $c$ , and that  $(E, \sqsubseteq_E)$  is a join semilattice.

## 5.2 Finding function $K$

In this subsection we derive a definition of function  $K$  satisfying (17), i.e., we construct a function  $K$  such that

$$agn\ (n+1) = K\ (agn\ n)$$

In the next subsection we show that  $K$  is monotonic with respect to  $\sqsubseteq_E$ , which allows us to conclude that  $agn\ \infty$  is the least fixed point of function  $K$ , i.e.,

$$agn\ \infty = \mu K$$

We calculate as follows for  $agn\ (n+1)$ . The goal is to express  $agn\ (n+1)$  in terms of  $agn\ n$ .

$$\begin{aligned} & agn\ (n+1) \\ = & \text{definition of } agn \\ & (id \triangle (af \cdot gh\ (n+1) \cdot N)) * (nts\ g) \end{aligned}$$

We proceed the calculation with the subexpression  $af \cdot gh\ (n+1) \cdot N$ . Suppose we can find a function  $J$  such that

$$af \cdot gh\ (n+1) \cdot N = J\ (agn\ n) \quad (18)$$

then we have

$$agn\ (n+1) = (id \triangle (J\ (agn\ n))) * (nts\ g)$$

and it follows by abstracting from  $agn\ n$  in the right-hand side of this equation that a function  $K$  satisfying (17) is defined by

$$K\ x = (id \triangle J\ x) * (nts\ g) \quad (19)$$

It remains to find a function  $J$  such that equation (18) is satisfied.

Function  $J$  satisfying equation (18) is obtained by manipulating the expression  $af\ (gh\ (n+1)\ (N\ nt))$ , where  $nt$  is an element of  $nts\ g$ . Abbreviate the right-hand argument of  $\oplus$  in definition (15) of  $gh\ (n+1)$  to  $rh$ , that is

$$\begin{aligned} rh = & [Node\ nt\ c \mid rhs \leftarrow rhss\ g\ nt, \\ & c \leftarrow cp\ ((gh\ n) * rhs)] \end{aligned}$$

Using this abbreviation we calculate as follows for  $af\ (gh\ (n+1)\ (N\ nt))$ .

$$\begin{aligned} & af\ (gh\ (n+1)\ (N\ nt)) \\ = & \text{definition of } gh \\ & af\ (gh\ n\ (N\ nt) \oplus rh) \\ = & \text{definition of } af \\ & (combine \cdot properties)\ (gh\ n\ (N\ nt) \oplus rh) \\ = & \text{definition of } combine \text{ and } properties \\ & (\oplus / \cdot p*)\ (gh\ n\ (N\ nt) \oplus rh) \\ = & (6) \\ & \oplus / (p* (gh\ n\ (N\ nt)) \oplus p* rh) \\ = & (11) \\ & (\oplus / \cdot p*)\ (gh\ n\ (N\ nt)) \oplus (\oplus / \cdot p*)\ rh \\ = & af, combine, \text{ and } properties \\ & af\ (gh\ n\ (N\ nt)) \oplus af\ rh \end{aligned}$$

We express the arguments of operator  $\oplus$  in the last expression above in terms of  $agn\ n$  separately. For  $agn\ n$  we have

$$agn\ n = (id \triangle (af \cdot gh\ n \cdot N)) * (nts\ g)$$

```

ag :: (Eq a
      ,Eq [(a,c)]
      ,Eq [c], Semilattice c
      ) =>
Grammar a b ->
(a -> [(Symbol a b,c)] -> c) ->
(b -> c) ->
[(a,c)]

```

## 5 The derivation of an algorithm

Function  $ag$  can be implemented in a functional language, but executing  $ag\ g\ p \oplus$  will result in a nonterminating computation because of the occurrence of  $\infty$  in the definition of function  $ag$ . This section derives an algorithm that can be implemented as an always terminating program that returns the value of  $ag\ g\ p \oplus$ . To obtain this algorithm we use the lattice theory given in Section 3.

Function  $ag$  satisfies the following equality.

$$\begin{aligned}
& ag\ g\ p \oplus \\
= & \\
& agn\ \infty \\
\text{where} & \\
& agn\ n = (id \triangle (af \cdot gh\ n \cdot N))^* (nts\ g) \\
& af = combine \cdot properties \\
& properties = p^* \\
& combine = \oplus /
\end{aligned}$$

This expression is obtained from the definition of function  $ag$  in Definition (16) by replacing the constant  $\infty$  by a variable  $n$ .

The CPO fixed point theorems may be used to find the value of  $agn\ \infty$  in finite time. Suppose there exists a function  $K$  such that for  $n \geq 0$

$$agn\ (n+1) = K\ (agn\ n) \quad (17)$$

If we suppose furthermore that there exists a CPO  $(E, \sqsubseteq_E)$  with bottom  $agn\ 0$ , then the results in Section 3 show that function  $K : E \rightarrow E$  has a least fixed point  $\mu K$ , defined by

$$\mu K = \sqcup / [K^n\ (agn\ 0) \mid n \leftarrow [0..]]$$

provided function  $K$  is continuous, and

$$agn\ \infty = \mu K$$

The domains used in the grammar analysis problems are finite, that is, the target type  $E$  of function  $ag$  is a finite type. Since every finite join semilattice is a CPO, and since each monotonic function on a finite domain is continuous, it suffices to find a join semilattice with bottom  $agn\ 0$ , and a monotonic function  $K$  satisfying (17).

This section consists of five subsections. The first subsection constructs a join semilattice with bottom  $agn\ 0$  for bottom-up grammar analysis problems. The second subsection derives a definition of function  $K$  that satisfies equation (17), i.e., it expresses  $agn\ (n+1)$  in terms of  $agn\ n$ . In order to find a definition of function  $K$  that satisfies equation (17) it will be advantageous to impose conditions upon the components of the bottom-up grammar analysis problem. The third subsection shows that provided some further conditions are satisfied the function  $K$  obtained in the derivation is monotonic. The fourth subsection discusses the conditions imposed thus far, and the fifth subsection applies the derived theory to some examples.

### 5.1 Constructing a join semilattice with bottom $agn\ 0$

We want to construct a join semilattice  $(E, \sqsubseteq_E)$  with bottom  $agn\ 0$  and join  $\sqcup_E$ , such that there exists a monotonic function  $K : E \rightarrow E$  satisfying  $agn\ (n+1) = K\ (agn\ n)$ . For that purpose, we impose our first condition on bottom-up grammar analysis problems.

For value  $agn\ 0$  we calculate as follows.

$$\begin{aligned}
& agn\ 0 \\
= & \text{definition of } agn \\
& (id \triangle (af \cdot gh\ 0 \cdot N))^* (nts\ g) \\
= & \text{definition of } gh, a^\bullet b = a \text{ for all } a \\
& (id \triangle (af \cdot [ ]^\bullet))^* (nts\ g) \\
= & f \cdot a^\bullet = (f\ a)^\bullet \\
& (id \triangle (af\ [ ]^\bullet))^* (nts\ g) \\
= & \text{definition of } af \\
& (id \triangle 1_{\oplus}^\bullet)^* (nts\ g)
\end{aligned}$$

We have derived the following equality.

$$agn\ 0 = (id \triangle 1_{\oplus}^\bullet)^* (nts\ g)$$

i.e.,  $agn\ 0$  is a list of length equal to the number of nonterminals of  $g$ , of which the second components

determine whether or not it is possible to derive the empty string from nonterminal  $nt$ , we apply the function  $combine$  to the list of results obtained by applying function  $p$  to all derivation trees with  $nt$  in the root. Function  $combine$  is defined by

$$combine = \vee/$$

Note that function  $combine$  corresponds with function  $([] \neq)$  occurring in the definition of  $Empty\ g\ nt$ , i.e., we have

$$\begin{aligned} [] \neq [] &= false \\ ([] \neq) (x ++ y) &= ([] \neq x) \vee ([] \neq y) \end{aligned}$$

Generalising this pattern, we now define the class of bottom-up grammar analysis problems.

**(16) Definition** *A bottom-up grammar analysis problem, which analyses a grammar  $g$  with respect to a function  $p : Rosetree\ a\ b \rightarrow c$ , and an operator  $\oplus : c \times c \rightarrow c$  with unit  $1_{\oplus}$ , is an expression of the form  $ag\ g\ p\ \oplus$ , where function  $ag$  is defined as follows.*

$$\begin{aligned} ag\ g\ p\ \oplus &= \\ &(id \triangle (af \cdot gh \infty \cdot N)) * (nts\ g) \\ \text{where} & \\ af &= combine \cdot properties \\ properties &= p * \\ combine &= \oplus/ \end{aligned}$$

In case of the bottom-up grammar analysis problem `EMPTY` we may now write

$$empties\ g = ag\ g\ (([] =) \cdot sen) \vee$$

and in case of the bottom-up grammar analysis problem `FIRST` we write

$$firsts\ g = ag\ g\ (take\ 1 \cdot sen) \cup$$

where operator  $\cup$  is defined by  $x \cup y = nub\ (x ++ y)$ .

## 4.4 Implementation

The definitions given above are translated into Gofer as follows.

We do not give the implementation of function  $ag$ , because the resulting program would not terminate. In the following subsection we transform function  $ag$  such that it always terminates, and

the resulting program can be found at the end of the following section.

We give the definitions of functions  $empties$  and  $firsts$ . There are some differences with the definitions given above. First, the last argument of the functions given above does not appear in the Gofer definitions below. We will assume later that the last argument of a bottom-up problem is always the join of a semi-lattice, so we need not pass it as an argument. The second difference is that function  $p$  is replaced by two functions, which are obtained by writing  $p$  as something very much like a catamorphism on  $Rosetree$ , namely, we will assume later that there exist functions  $pn$  and  $pl$  such that

$$\begin{aligned} p\ (Node\ nt\ x) &= pn\ nt\ ((top \triangle p) * x) \\ p\ (Leaf\ x) &= pl\ x \end{aligned}$$

The second argument of  $ag$  is derived from function  $pn$ , and the third argument of  $ag$  is the function  $pl$ . At the end of the following section we construct these definitions.

```
empties :: Eq [(a,Bool)] =>
  Grammar a b ->
  [(a,Bool)]

empties g =
  ag
  g
  (\nt x -> and (map snd x))
  (\a -> False)

firsts :: (Eq [(a,Bool)]
  ,Eq [(a,[b])])
  ,Eq [[b]]
  ,Semilattice [b]
  ) =>
  Grammar a b ->
  [(a,[b])]

firsts g =
  ag
  g
  (\nt x -> foldr t bottom x)
  (\b -> [b])
  where
  t (N a,y) x
    | eg 'at' a = nub (y ++ x)
    | otherwise = y
  t (T b,y) x = y
  eg = empties g
```

From these definitions we obtain the following type for function  $ag$ .

Before we give the definition, we discuss the function  $cp$  (cartesian product), which is used in the definition of function  $generate$ .

**Function  $cp$**

Function  $cp$  returns the cartesian product of a list of lists. It is defined as a map followed by a reduce by

$$\begin{aligned} cp &= \lambda/ \cdot [ ]^{**} \\ xs \lambda ys &= [x \# y \mid x \leftarrow xs, y \leftarrow ys] \end{aligned}$$

where  $[ \cdot ]$  takes an element  $a$ , and returns the singleton list containing that element:  $[a]$ . Note that  $[ [ ] ]$  is the unit of operator  $\lambda$ . Function  $cp$  commutes with function  $f^{**}$  for all functions  $f$ , i.e., for all functions  $f$  we have

$$f^{**} \cdot cp = cp \cdot f^{**} \quad (14)$$

**Function  $generate$**

Function  $gh$  is defined in the context of a grammar  $g$ , which from now on is considered a constant. It takes a natural number  $n$ , and a symbol  $nt$ , and returns the collection of all derivation trees, of height at most  $n$ , derivable with the productions of  $g$  with symbol  $nt$  in the root, so

$$\begin{aligned} gh : nat \rightarrow Symbol \ a \ b \rightarrow (Rosetree \ a \ b)^* \\ gh \ n \ nt &= [y \mid nt \xrightarrow{*} y \wedge height \ y \leq n] \end{aligned}$$

where we suppose that  $\xrightarrow{*}$  derives derivation trees instead of strings with productions from grammar  $g$ . So  $dt$  is an element of  $gh \ 5 \ (N \ E)$ . Function  $generate$  is defined in terms of function  $gh$  as follows.

$$generate \ g = (gh \ \infty \cdot N)^* \ (nts \ g)$$

Function  $gh$  can be defined recursively in various ways; we have chosen the following definition which is easily manipulated in calculations. Function  $gh$  is defined by pattern matching on its first argument, using the second argument to break the tie. There are no trees of height zero, so

$$gh \ 0 \ symbol = [ ]$$

There is just one derivation tree of height at most  $n+1$  that can be built from a terminal.

$$gh \ (n+1) \ (T \ b) = [Leaf \ b]$$

The list of derivation trees of height at most  $n+1$  derivable from a nonterminal  $nt$  contains the list of the derivation trees of height at most  $n$  derivable from  $nt$ . Furthermore, for each production

for  $nt$  we add the cartesian product of the derivation trees of height at most  $n$  of the symbols of the right-hand side of a production for  $nt$ ; each element of the cartesian product is turned into a derivation tree using function  $Node \ nt$ .

$$\begin{aligned} gh \ (n+1) \ (N \ a) \\ = \\ (gh \ n \ (N \ a)) \# [Node \ a \ c \mid \\ rhs \leftarrow rhss \ g \ a, c \leftarrow cp \ ((gh \ n)^* \ rhs)] \end{aligned} \quad (15)$$

We do not bother about duplicate elements in  $gh \ n \ nt$ ; applying function  $nub$  to the right-hand expression of the last equation would have removed them. The right-hand side argument of  $\#$  in the last equation of the definition of function  $gh$  can be rewritten using laws for list-comprehensions.

$$\begin{aligned} [Node \ nt \ c \mid \\ rhs \leftarrow rhss \ g \ nt, c \leftarrow cp \ ((gh \ n)^* \ rhs)] \\ = \text{equation (9) for list-comprehensions} \\ (Node \ nt)^* \\ [c \mid rhs \leftarrow rhss \ g \ nt, c \leftarrow cp \ ((gh \ n)^* \ rhs)] \\ = \text{(10), (8), and (9)} \\ ((Node \ nt)^* \cdot \# / \cdot cp^* \cdot (gh \ n)^{**}) \\ [rhs \mid rhs \leftarrow rhss \ g \ nt] \\ = \text{(8), definition of function } rhss \\ ((Node \ nt)^* \cdot \# / \cdot cp^* \cdot (gh \ n)^{**}) \\ (rhss \ g \ nt) \end{aligned}$$

This equality will be used in the calculation in Section 5.

### 4.3 Bottom-up problems

We formalise the notion of a grammar analysis problem. In case of the EMPTY problem, we want to determine for all nonterminals  $nt$  from a grammar  $g$  whether or not it is possible to derive the empty string from nonterminal  $nt$ . A non-executable specification for this problem reads as follows. Given a nonterminal  $nt$  we apply a function  $p$  to each derivation tree with  $nt$  in the root. Function  $p$  determines whether or not the string represented by the derivation tree is empty, i.e.,

$$p = ([ ] =) \cdot sen$$

Note that function  $p$  corresponds with the two expression  $nt \xrightarrow{*} x, x = [ ]$  occurring in the list-comprehension in the definition of  $Empty \ g \ nt$ . To

## 4 Grammar analysis problems

Although in some grammar analysis problems only a property of the start-symbol of the grammar is sought, we define a grammar analysis problem to be a problem which requires finding information about all nonterminals of the grammar. This section defines bottom-up grammar analysis problems. The first subsection gives some examples of grammar analysis problems. The second subsection discusses functions for generating derivation trees. The third subsection defines bottom-up grammar analysis problems, and, finally, the fourth subsection gives the implementation of some of the functions introduced in this section.

### 4.1 Examples of grammar analysis problems

Part of determining whether or not a grammar is LL(1) consists of computing lookahead sets. If the grammar analysis problems `EMPTY`, `FIRST`, and `FOLLOW` have been solved, we can easily approximate the lookahead sets. The definitions of these problems are our first three examples. The fourth example, `LEFT-CONTEXT` concerns the computation of left-contexts of nonterminals. Left-contexts of nonterminals are used to determine whether or not a grammar is LR(0).

#### `EMPTY`

Given a grammar  $g$  and a nonterminal  $nt$  from  $g$ , the expression *Empty  $g$   $nt$*  is a boolean expression whether or not it is possible to derive the empty string from  $nt$ , using the productions from  $g$ . It is defined by

$$\begin{aligned} & \textit{Empty } g \textit{ } nt \\ = & \quad [ ] \neq [x \mid nt \xrightarrow{*} x, x = [ ] ] \end{aligned}$$

where  $\xrightarrow{*}$  denotes a derivation with productions from  $g$ .

#### `FIRST`

Given a grammar  $g$  and a nonterminal  $nt$  from  $g$ , the expression *First  $g$   $nt$*  is the set of terminals that can appear as the first element of a sentence derivable from  $nt$ . It is defined by

$$\begin{aligned} & \textit{First } g \textit{ } nt \\ = & \quad \textit{nub } [a \mid nt \xrightarrow{*} [a] \textit{ } \textit{ } x, x \in X^*] \end{aligned}$$

where  $X$  is the set of terminals of  $g$ .

#### `FOLLOW`

Given a grammar  $g$  and a nonterminal  $nt$  from  $g$ , the expression *Follow  $g$   $nt$*  is the set of terminals that can follow on  $nt$  in a derivation starting with the start-nt  $S$  from  $g$ . It is defined by

$$\begin{aligned} & \textit{Follow } g \textit{ } nt \\ = & \quad \textit{nub } [a \in X \mid S \xrightarrow{*} u \textit{ } \textit{ } [nt, a] \textit{ } \textit{ } v] \end{aligned}$$

#### `LEFT-CONTEXT`

A left-context of a nonterminal is a list of terminals and nonterminals that can appear before the nonterminal in a right-most derivation from the start-nt, provided the list of symbols after the nonterminal is a list of terminals. Given a grammar  $g$  and a nonterminal  $nt$  from  $g$ , the expression *LC  $g$   $nt$*  is the set of left-contexts of  $nt$ .

$$\begin{aligned} & \textit{LC } g \textit{ } nt \\ = & \quad \textit{nub } [u \mid S \xrightarrow{*} u \textit{ } \textit{ } [nt] \textit{ } \textit{ } v, v \in X^*] \end{aligned}$$

### Bottom-up versus top-down

The definitions in the first two examples given above require finding information about a nonterminal, and do not refer to the context in which such a nonterminal appears. These two examples are bottom-up grammar analysis problems. The definitions of the last two examples explicitly refer to the context in which the nonterminal appears, namely  $u \textit{ } \textit{ } [_, a] \textit{ } \textit{ } v$ , and  $u \textit{ } \textit{ } [_, ] \textit{ } \textit{ } v$ , respectively. These two examples are top-down grammar analysis problems. In the rest of the paper we limit ourselves to bottom-up problems.

## 4.2 Generating trees

The definitions in the examples of grammar analysis problems given in the previous subsection typically refer somehow to all sentences derivable from a nonterminal. The sentences derivable from a nonterminal can be obtained from the derivation trees of the grammar with the given nonterminal in the root. In this subsection we define a function returning all possible derivation trees of a grammar.

Function *generate* takes a grammar, and returns a list of lists, in which each list contains all derivation trees with the same nonterminal in the root.

Note that the join of two elements in  $D$  is uniquely defined when it exists. The *least upperbound* or *lub* of a subset  $X \subseteq D$  is denoted by  $\sqcup / X$ . It is defined by

$$\begin{aligned} c &= \sqcup / X \\ \equiv & (\forall e :: c \sqsubseteq e \equiv (\forall x : x \in X : x \sqsubseteq e)) \end{aligned}$$

Not every  $X \subseteq D$  needs to have a lub. The *meet*  $\sqcap$  and *greatest lowerbound* or *glb* are dual to the join and the lub, respectively. Their definitions are omitted. Let  $(D, \sqsubseteq)$  be a poset. If for all elements  $d$  and  $d'$  their join  $d \sqcup d'$  exists, then  $(D, \sqsubseteq)$  is called a *join semilattice*. A *meet semilattice* is defined similarly. Let  $S$  be a subset of a poset.  $S$  is said to be *directed* if every finite subset of  $S$  has an upper bound. A poset  $D$  is a *complete partial order* or *cpo* if it contains a bottom element, and if each directed subset of  $D$  has a lub, so  $\sqcup / X$  exists for all directed subsets  $X \subseteq D$ .

### Fixed Points

An element  $d \in D$  is a *fixed point* of function  $f : D \rightarrow D$  if  $f d = d$ . It is a *least fixed point* if for any other fixed point  $d'$  of  $f$  we have  $d \sqsubseteq d'$ . A function  $f : D \rightarrow E$  is *monotonic* if it respects the ordering on  $D$ , i.e.,

$$d \sqsubseteq d' \Rightarrow f d \sqsubseteq f d'$$

A function  $f : D \rightarrow E$  is *continuous* if it respects lubs of directed subsets, i.e., if  $X \subseteq D$  is a directed subset, then

$$(f \cdot \sqcup / X) = (\sqcup / f^* X)$$

Let  $D$  be a finite set,  $(D, \sqsubseteq)$  a CPO with bottom  $\perp$ , and  $g : D \rightarrow D$  a continuous function. It follows from the CPO Fixed Point Theorem I [2] that function  $g$  has a least fixed point  $\mu g$ , defined by

$$\mu g = \sqcup / [g^n \perp \mid n \leftarrow [0..]]$$

Since  $g^i \perp \sqsubseteq g^{i+1} \perp$ , we have that the least fixed point of  $g$  equals the first element in  $[g^n \perp \mid n \leftarrow [0..]]$  that occurs twice, i.e.,

$$\sqcup / [g^n \perp \mid n \leftarrow [0..]] = \text{lf}p g \perp$$

where function  $\text{lf}p$  is defined by

$$\text{lf}p f x = \begin{cases} x & \text{if } f x = x \\ \text{lf}p f (f x) & \text{otherwise} \end{cases}$$

The *Fixed Point Fusion* Theorem (or Plotkin's Lemma) is used to reason about fixed points. This theorem reads as follows.

$$\begin{aligned} f \perp &= \perp \wedge f \cdot h = g \cdot f \\ \Rightarrow & \\ f \mu h &= \mu g \end{aligned}$$

We use the Fixed Point Fusion Theorem and the CPO Fixed Point Theorem I as follows. Consider the function  $(+1)$ . Define  $\infty = \mu(+1)$ . Taking  $h = (+1)$ , applying the Fixed Point Fusion Theorem gives

$$\begin{aligned} f \perp &= \perp \wedge f \cdot (+1) = g \cdot f \\ \Rightarrow & \\ f \mu(+1) &= \mu g \end{aligned}$$

Writing  $\infty$  for  $\mu(+1)$ , and 0 for the bottom  $\perp$  of the natural numbers, we get

$$\begin{aligned} f 0 &= \perp \wedge f (n+1) = g (f n) \\ \Rightarrow & \\ f \infty &= \mu g \end{aligned}$$

## 3.1 Implementation

The definitions of some of the functions and classes given above are translated into Gofer as follows.

```
class Semilattice a where
  join    :: a -> a -> a
  bottom :: a

instance Semilattice Bool where
  join = (||)
  bottom = False

instance Eq a => Semilattice [a] where
  join = \a b -> nub (a ++ b)
  bottom = []

lf p :: Eq a =>
  (a -> a) ->
  a ->
  a
lf p f x | x == f x = x
          | otherwise = lf p f (f x)

lub :: Semilattice a => [a] -> a
lub = foldl join bottom
```

corresponds to the derivation tree  $dt$  defined by

```

Node E
  [Node E
    [Node T
      [Node F
        [Leaf v]]],
    Leaf +,
    Node T
      [Node F
        [Leaf v]]
  ]

```

Suppose function  $top : Rosetree\ a\ b \rightarrow Symbol\ a\ b$  returns the top of a rose-tree. For each subtree of a derivation tree of the form  $Node\ a\ x$  we have that  $a \rightarrow top\ x$  is a production of the grammar.

The function  $sen$  takes a rose-tree, and returns the sentence of which the rose-tree is a derivation. Function  $sen$  is defined by

```

sen (Node a x) = ++/ (sen* x)
sen (Leaf b)   = [b]

```

It follows that  $sen\ dt = v+v$ .

### Catamorphisms on Rose-Trees

For every recursive datatype we can define a function which recursively replaces constructors by functions [6]. By definition, a *catamorphism* on the datatype  $Rosetree$  is a function  $h : Rosetree\ a\ b \rightarrow c$  that is uniquely determined by functions  $f$  and  $g$  as follows.

```

h (Node a x) = f a (h* x)
h (Leaf b)   = g b

```

For such a function  $h$  we write

```

h = RT_cata f g

```

The function  $sen$  defined above is a catamorphism, i.e.,

```

sen = RT_cata s t
  where s a x = ++/ x
        t b = [b]

```

Another example of a catamorphism on  $Rosetree$  is the *height* function, which returns the height of a rose-tree.

```

height = RT_cata s t
  where s (N a) x = 1 + ↑/ x
        t (T b) = 1

```

where  $\uparrow$  returns the maximum of two numbers. For example,  $height\ dt = 5$ . We will encounter several other catamorphisms on rose-trees in the following sections.

## 2.1 Implementation

The definitions of some of the functions and datatypes given above are translated into Gofer as follows.

```

split f g x = (f x , g x)

data Symbol a b = N a
                | T b

type Grammar a b = (a, [(a, [Symbol a b])])

rhss :: Eq a =>
  Grammar a b ->
  a ->
  [[Symbol a b]]
rhss g nt = [rhs
             | (z,rhs) <- snd g
             , z==nt
             ]

```

```

nts :: Eq a => Grammar a b -> [a]
nts g = nub (map fst (snd g))

```

## 3 Lattice theory

This section only gives the definitions of notions from lattice theory that are used in the subsequent sections. For a more extensive introduction to lattice theory the reader is referred to e.g. [2].

### Lattices and CPO's

A *partial order* on a set  $A$  is a reflexive, antisymmetrical, and transitive binary relation on  $A$ . A *partially ordered set* or *poset* is a pair  $(D, \sqsubseteq)$  consisting of a set  $D$  together with a partial order  $\sqsubseteq$  on  $D$ . If it exists, the least or *bottom* element of a poset is usually denoted by  $\perp$ . Given  $d, d' \in D$ , their *join*, denoted by  $d \sqcup d'$ , is the least element in  $D$  that is greater than both  $d$  and  $d'$ . It is fully characterised by the following equation:

$$c = d \sqcup d' \equiv (\forall e :: c \sqsubseteq e \equiv d \sqsubseteq e \wedge d' \sqsubseteq e)$$

Function *concat* is an example of a reduce: it is defined by

$$\text{concat} = \text{+}/$$

For the composition of a map and function *+/*, and for the composition of a reduce and function *+/* we have

$$f* \cdot \text{+}/ = \text{+}/ \cdot f** \quad (12)$$

$$\oplus/ \cdot \text{+}/ = \oplus/ \cdot (\oplus/)* \quad (13)$$

### Terminals and nonterminals

Suppose a terminal is a value of type *b*, and a nonterminal is a value of type *a*. A symbol that is either a nonterminal of type *a* or a terminal of type *b* is a value of the datatype *Symbol* defined as

$$\text{data Symbol } a \ b = N \ a \mid T \ b$$

An element *N x* is considered to be a nonterminal, and an element *T y* is considered to be a terminal. Note that this definition makes the conventional disjoint sum operational.

### Grammars

A context-free grammar consists of sets of nonterminals, terminals, productions, and a start-symbol. In Gofer, the sets of nonterminals and terminals correspond with the types *a* and *b*, respectively. These types are parameters of the definition of a context-free grammar. We represent a context-free grammar in Gofer by a pair, the first component of which denotes the start-symbol, and the second component of which denotes the productions of the grammar. The start-symbol is a nonterminal, i.e., a value of type *a*. The productions of a grammar are a set (represented by a list) of pairs the left-component of which is a nonterminal, and the right component of which is a list of symbols. A context-free grammar is a value of the type *Grammar*, which is defined by

$$\begin{aligned} & \text{type Grammar } a \ b \\ = & (a, [(a, [Symbol \ a \ b])]) \end{aligned}$$

For example, consider the grammar *eg* for expressions in a variable *v*.

$$\begin{aligned} E & \rightarrow E+T \mid T \\ T & \rightarrow T*F \mid F \\ F & \rightarrow (E) \mid v \end{aligned}$$

This grammar is represented in Gofer by

$$\begin{aligned} eg = & (E, [(E, [N \ E, T \ +, N \ T]) \\ & , (E, [N \ T]) \\ & , (T, [N \ T, T \ *, N \ F]) \\ & , (T, [N \ F]) \\ & , (F, [T \ (, N \ E, T \ )]) \\ & , (F, [T \ v])]) \\ & ) \end{aligned}$$

Function *rhss* takes a grammar and a nonterminal *nt* and returns the right-hand sides of the productions of *nt*. It is defined by

$$\text{rhss } g \ nt = [rhs \mid (nt, rhs) \leftarrow \text{exr } g]$$

Function *nts* takes a grammar, and returns the list of nonterminals of the grammar. We assume that for each nonterminal there exists at least one production. Function *nts* is defined by

$$\text{nts } g = \text{nub } (\text{exl}* (\text{exr } g))$$

where function *nub* removes duplicates from a list.

### Parse Trees

To determine whether or not the empty string can be derived from a nonterminal (the EMPTY problem), we have to refer to all sentences that are derivable from the given nonterminal in the given grammar. A derivation using productions of a context-free grammar corresponds to a *parse tree* or *derivation tree*, i.e., an element of the datatype *Rosetree*, which is defined by

$$\begin{aligned} & \text{data Rosetree } a \ b \\ = & \text{Node } a \ [Rosetree \ a \ b] \mid \text{Leaf } b \end{aligned}$$

For example, the following derivation of the sentence *v+v* using the productions from grammar *eg*

$$\begin{aligned} & E \\ \Rightarrow & E+T \\ \Rightarrow & T+T \\ \Rightarrow & F+T \\ \Rightarrow & v+T \\ \Rightarrow & v+F \\ \Rightarrow & v+v \end{aligned}$$

This paper is organised as follows. Section 2 defines the datatypes that are used in manipulating grammars in Gofer. Section 3 introduces some necessary concepts of lattice theory needed in the subsequent sections. Section 4 defines the class of bottom-up grammar analysis problems, and gives some examples. Section 5 derives an algorithm that can be used to solve bottom-up problems. Section 6 concludes the paper.

## 2 Datatypes for grammars in Gofer

This section defines various datatypes in Gofer used in analysing and representing grammars.

### Functions

We use simple juxtaposition and a little white space to denote the application of a function  $f : s \rightarrow t$  to an argument  $x \in s$ , i.e.,  $f x$ . Composition of functions  $f : s \rightarrow t$  and  $g : r \rightarrow s$  is written  $f \cdot g : r \rightarrow t$ . Composition is associative, and the identity function  $id$  is the unit of composition. Projection  $exl$  ( $exr$ ) selects the left (right) component of a pair, i.e.,

$$\begin{aligned} exl(a, b) &= a \\ exr(a, b) &= b \end{aligned}$$

Given functions  $f : A \rightarrow B$  and  $g : A \rightarrow C$ , function  $f \triangle g : A \rightarrow B \times C$  (split) applies both  $f$  and  $g$  to an argument. The type  $B \times C$  is the cartesian product of the sets  $B$  and  $C$ .

$$(f \triangle g) a = (f a, g a)$$

Given functions  $f : A \rightarrow B$  and  $g : C \rightarrow D$ , function  $f \times g : A \times C \rightarrow B \times D$  (product) applies  $f$  to the first component, and  $g$  to the second component of its argument.

$$(f \times g) (a, c) = (f a, g c)$$

We have the following laws concerning projections, split, and product.

$$exl \cdot (f \triangle g) = f \tag{1}$$

$$exr \cdot (f \triangle g) = g \tag{2}$$

$$exr \cdot f \times g = g \cdot exr \tag{3}$$

$$exl \cdot f \times g = f \cdot exl \tag{4}$$

$$f \times g \cdot h \triangle j = (f \cdot h) \triangle (g \cdot j) \tag{5}$$

Function application binds stronger than a binary operator, and among the binary operators function composition binds weakest.

### Lists

The datatype *list* is a prominent datatype in the subsequent sections, and we will use a number of properties that are satisfied by functions defined on the datatype *list*. The empty list is denoted by  $[]$ , and the concatenation of two lists  $x$  and  $y$  is denoted by  $x \# y$ . Prepending an element  $x$  to a list  $xs$  is denoted by  $x : xs$ . The datatype *list* over base type  $A$  is denoted by  $A^*$ . For  $f : A \rightarrow B$ , function  $f^* : A^* \rightarrow B^*$ , called a *map* function takes a list and applies function  $f$  to all elements in the list, so

$$f^* xs = [f x \mid x \leftarrow xs]$$

For the map function we have

$$\begin{aligned} f^* [] &= [] \\ f^* (x \# y) &= f^* x \# f^* y \\ f^* (x : xs) &= f x : f^* xs \end{aligned} \tag{6}$$

Map-distributivity says that the composition of two maps is a map again, i.e., for all functions  $f$  and  $g$ :

$$f^* \cdot g^* = (f \cdot g)^* \tag{7}$$

Furthermore, the result of mapping the identity function over an argument is the argument itself, so

$$id_{A^*} = id_{A^*}$$

These equalities say that  $*$  is a *functor*. An important functional programming construct we use is *list-comprehension*. For example,

$$\begin{aligned} &[(x, y) \mid x \leftarrow [1, 2], y \leftarrow [3, 4]] \\ &= \\ &[(1, 3), (1, 4), (2, 3), (2, 4)] \end{aligned}$$

We will use the following laws for list-comprehensions [11] in some calculations.

$$[t \mid t \leftarrow ts] = ts \tag{8}$$

$$[f t \mid q] = f^* [t \mid q] \tag{9}$$

$$[t \mid p, q] = concat [[t \mid q] \mid p] \tag{10}$$

where function *concat* flattens a list of lists. Function *concat* is defined in terms of the reduce operator. The *reduce* operator  $/$  takes an associative operator  $\oplus$  with unit  $1_\oplus$ , and a list, and places the operator in between the elements of a list, so  $\oplus / [a, b, c] = a \oplus b \oplus c$ . For operator  $\oplus : A \times A \rightarrow A$  we have  $\oplus / : A^* \rightarrow A$ . It is defined by

$$\begin{aligned} \oplus / [] &= 1_\oplus \\ \oplus / (x \# y) &= \oplus / x \oplus \oplus / y \\ \oplus / (x : xs) &= x \oplus \oplus / xs \end{aligned} \tag{11}$$

# Bottom-up Grammar Analysis

## — A Functional Formulation —\*

Johan Jeuring and Doaitse Swierstra

Utrecht University  
P.O.Box 80.089, 3508 TB Utrecht, The Netherlands  
Email: {johan,doaitse}@cs.ruu.nl

### Abstract

This paper discusses bottom-up grammar analysis problems such as the EMPTY problem and the FIRST problem. It defines a general class of bottom-up grammar analysis problems, and from this definition it derives a functional program for performing bottom-up grammar analysis. The derivation is purely calculational, using theorems from lattice theory, the Bird-Meertens calculus, and laws for list-comprehensions. Sufficient conditions guaranteeing the existence of a solution emerge as a byproduct of the calculation. The resulting program is used to construct programs for the EMPTY problem and the FIRST problem.

### 1 Introduction

Grammar analysis is performed in many different situations: Yacc tests whether or not its input grammar is LALR(1), parser generators contain functions for determining whether or not a nonterminal can derive the empty string (EMPTY) as part of determining the set of all symbols that can appear as the first symbol of a derived string (FIRST), and for determining the set of symbols that can appear as the first symbol following upon a string derived by a given nonterminal (FOLLOW). Other, similar, problems arise when analysing attribute dependencies in attribute grammars: determine the inherited attributes upon which a synthesised attribute depends (IS), and, conversely, determine the synthesised attributes upon which an inherited attribute depends (SI). Such problems are called *grammar analysis problems*. Grammar analysis problems can be divided

into two classes: *bottom-up* and *top-down*. The difference between these classes is that the required information for a nonterminal in a top-down problem depends on the possible contexts for that nonterminal, whereas in a bottom-up problem the contexts of a nonterminal can be ignored. Often the output of a bottom-up problem is used in a top-down problem. The specification of a grammar analysis problem determines the class to which it belongs: EMPTY, FIRST, and IS are bottom-up grammar analysis problems, the FOLLOW and SI problems belong to the top-down class. This paper studies bottom-up grammar analysis.

Grammar analysis problems are described by sets of mutually recursive equations, and the solution of a grammar analysis problem is a fixed point of this equational system. Möncke and Wilhelm [9] observe this, and give several solutions, depending on the conditions that are satisfied, for such problems. The goal of this paper is to *derive* the solutions given by Möncke and Wilhelm. We start with a very general specification of a bottom-up grammar analysis problem, and we derive a function of which the fixed point gives the solution of the problem. This function is obtained by applying laws to components of the expressions occurring in the specification. The laws we apply are familiar laws for, for example, list-comprehensions [11], and maps [1, 7]. Sufficient conditions for guaranteeing the existence of a fixed point solution emerge as a byproduct of this derivation. Finally we give the implementation of the derived algorithm in the functional language Gofer [5, 3]. Incorporating the functions for solving grammar analysis problems in parser generators such as a functional version of Yacc [10] and Ratatosk [8] would reduce the amount of code used in, and very likely increase the speed of, these parser generators.

---

\*This paper is an extended version of a paper with the same title that will be presented at ESOP '94.



# Bottom-up Grammar Analysis — A Functional Formulation —

Johan Jeuring and Doaitse Swierstra

Technical Report UU-CS-1994-01  
January 1994

Department of Computer Science  
Utrecht University  
P.O.Box 80.089  
3508 TB Utrecht  
The Netherlands

# Bottom-up Grammar Analysis — A Functional Formulation —

Johan Jeuring and Doaitse Swierstra

UU-CS-1994-01  
January 1994



**Utrecht University**

**Department of Computer Science**

Padualaan 14, P.O. Box 80.089,  
3508 TB Utrecht, The Netherlands,  
Tel. : + 31 - 30 - 531454