

A Boolean grammar for a simple programming language

Alexander Okhotin
okhotin@cs.queensu.ca

Technical report 2004-478
School of Computing, Queen's University,
Kingston, Ontario, Canada K7L 3N6

March 2004

Abstract

A toy procedural programming language is defined, and a Boolean grammar for the set of well-formed programs in this language is constructed. This is apparently the first specification of a programming language entirely by a formal grammar.

Contents

1	Introduction	2
2	Definition of the language	2
3	Boolean grammars	6
4	A Boolean grammar for the language	8
5	Conclusion	16
A	A parser	19
B	Test suite	22
C	Performance	33

1 Introduction

Formal methods have been employed in the specification of programming languages since the early days of computer science. Already the Algol 60 report [5, 6] used a context-free grammar to specify the complete lexical contents of the language, and, partially, its syntax. But still many actually syntactical requirements, such as the rules requiring the declaration of variables, were described in plain words together with the semantics of the language.

The hope of inventing a more sophisticated context-free grammar that would specify the entire syntax of Algol 60 was buried by Floyd [3], who used a formal language-theoretic argument based upon the non-context-freeness of the language $\{a^n b^n c^n | n \geq 0\}$ to prove that identifier checking cannot be performed by a context-free grammar, and hence the set of well-formed Algol 60 programs is not context-free.

The appealing prospects of giving a completely formal specification of the syntax of programming languages were not entirely abandoned, though. In the aftermath of Floyd's destructive result, many extensions of context-free grammars pertaining to cope with this problem were developed. One of the most promising attempts was represented by van Wijngaarden's two-level grammars [14, 15], which were actually used in the specification of Algol 68. Sadly, they were soon proved to be computationally universal, which demonstrated them hardly fit for practical use.

Among the formalisms proposed at that time that were more successful in evading Turing completeness, let us mention Aho's indexed grammars [1], which can denote all of the abstract non-context-free languages typically given to illustrate the non-context-freeness of programming languages [2]. However, it turned out that the possibility of denoting these constructs does not at all imply that a grammar for any programming language can in fact be constructed. Also, no polynomial-time parsing algorithm for indexed grammars has been found. As a result, indexed grammars remained an abstract model in the computation theory, while the syntax of new programming languages continued to be defined in the partially formal way introduced in the Algol 60 report [5, 6].

This paper gives a completely grammatical definition for the set of well-formed programs in a tiny procedural programming language. The recently introduced Boolean grammars [10] are used as the language specification formalism. The actual grammar is constructed and explained; this grammar has also been tested with a prototype parser generator [9], demonstrating square-time performance.

2 Definition of the language

Following is a complete definition of the syntax of the model language, i.e., a specification of the set of well-formed programs.

[I] Lexical conventions.

[I.1] A program is a finite string over the following alphabet of 55 characters:

[I.1.1] 26 letters: a, . . . , z.

[I.1.2] 10 digits: 0, . . . , 9.

[I.1.3] 2 whitespace characters: space `␣` and newline `¶`.

[I.1.4] 17 punctuators: (,), {, }, “, ”, “;”, +, -, *, /, %, &, |, !, =, <, >.

[I.2] The following finite strings over the alphabet [I.1] are called *lexemes*:

[I.2.1] Five keywords: `var`, `if`, `else`, `while`, `return`.

[I.2.2] An identifier: a finite nonempty sequence of letters [I.1.1] and digits [I.1.2] that starts from a digit and does not coincide with any of the keywords [I.2.1].

[I.2.3] A number: a finite nonempty sequence of digits [I.1.2].

[I.2.4] 13 binary infix operators: +, -, *, /, %, &, |, <, >, <=, >=, ==, !=.

[I.2.5] 2 unary prefix operators: -, !.

[I.2.6] Assignment operator: =.

[I.2.7] Parentheses: (,).

[I.2.8] Compound statement delimiters: {, }.

[I.2.9] Comma and semicolon: “,”, “;”.

[I.3] Separation of character stream into lexemes.

[I.3.1] The input [I.1] is processed from left to right.

[I.3.2] Each time the longest possible sequence of characters that forms a lexeme [I.2] is consumed, or a whitespace character [I.1.3] is discarded.

[I.3.3] Whitespace characters act as explicitly stated boundaries between the lexemes.

[I.3.4] According to the basic syntax defined below [II], the only essential uses of whitespace in a well-formed program are:

- in a declaration statement [II.2.3], to separate the keyword `var` [I.2.1] from the next identifier [I.2.2];
- in a conditional statement [II.2.4], to separate the keyword `else` [I.2.1] from a keyword, an identifier [I.2.2] or a number [I.2.3] that possibly start the inner statement [II.2];
- in a return statement [II.2.6], to separate the keyword `return` [I.2.1] from the next identifier [I.2.2] or a number [I.2.3] that possibly opens an expression [II.1].

[II] Basic syntax.

[II.1] Expressions.

- [II.1.1]** An identifier [I.2.2] is an expression.
- [II.1.2]** A number [I.2.3] is an expression.
- [II.1.3]** (e) is an expression for every expression e .
- [II.1.4]** Function call: $f (e_1 , \dots , e_k)$ is an expression for every identifier f [I.2.2] and expressions e_1, \dots, e_k ($k \geq 0$). This is a call of the function f with the arguments e_1, \dots, e_k .
- [II.1.5]** $e_1 \text{ op } e_2$ is an expression for every binary operator op [I.2.4] and expressions e_1, e_2 .
- [II.1.6]** $\text{op } e$ is an expression for every expression e and unary operator op [I.2.5].
- [II.1.7]** Assignment: $x = e$ is an expression for every identifier x [I.2.2] and expression e .

[II.2] Statements.

- [II.2.1]** Expression-statement: $e ;$ is a statement for every expression e [II.1].
- [II.2.2]** Compound statement: $\{ s_1 s_2 \dots s_k \}$ is a statement for every $k \geq 0$ and for all statements s_1, \dots, s_k .
- [II.2.3]** Declaration statement: $\text{var } x_1, \dots, x_k ;$ is a statement for every $k \geq 1$ and for all identifiers x_1, \dots, x_k [II.1.1].
- [II.2.4]** Conditional statement: $\text{if } (e) s$ and $\text{if } (e) s \text{ else } s'$ are statements for every expression e [II.1] and statements s, s' .
- [II.2.5]** Iteration statement: $\text{while } (e) s$ is a statement for every expression e [II.1] and statement s .
- [II.2.6]** Return statement: $\text{return } e ;$ is a statement for every expression e [II.1].

[II.3] Functions.

- [II.3.1]** Function header: $\mathbf{f} (x_1, \dots, x_k)$ is called a *function header* for every $k \geq 0$ and for all identifiers f, x_1, \dots, x_k [II.1.1]. The identifier f is the *name* of the function, the identifiers x_1, \dots, x_k are the *formal arguments* of the function.
- [II.3.2]** A *function declaration* is a function header [II.3.1] followed by a compound statement [II.2.2]. This compound statement is called the *body* of the function.

[II.4] A *program* is a finite sequence of function declarations [II.3.2].

[III] Additional syntactical constraints on a program [II.4].

[III.1] Function declarations.

- [III.1.1] A program may not contain multiple declarations of functions [II.3.2] sharing the same name [II.3.1].
- [III.1.2] A program must contain a declaration of a function [II.3.2] with the name `main`, which must have exactly one formal argument [II.3.1].
- [III.2] Calls to functions.
 - [III.2.1] Any call to a function f [II.1.4] should be preceded by a function header [II.3.1] that opens a declaration of a function with the name f [II.3.2].
 - [III.2.2] The number of arguments in a call to a function f [II.1.4] should match the number of formal arguments in the function header of f [II.3.1].
- [III.3] Variable declarations and their scope.
 - [III.3.1] *Variables* are declared in function headers [II.3.1] and in declaration statements [II.2.3]. Each identifier [I.2.2] in the corresponding lists declares a variable, and is called the *name* of the variable.
 - [III.3.2] With every variable declaration [III.3.1], a certain nonempty set of statements [II.2] is associated, called *the scope of the declaration*.
 - [III.3.3] The scope of a declaration of a variable as a formal argument of a function [II.3.1] is the whole body of the function [II.3.2].
 - [III.3.4] The scope of a declaration in a `var` statement [II.2.3] located within a compound statement [II.2.2] covers all statements in this compound statement starting from the `var` statement in question:

$$\{s_1 \dots s_{i-1} \underbrace{\text{var } x_1, \dots, x_k; s_{i+1} \dots s_k}_{\substack{\text{a declaration} \\ \text{the scope of this declaration}}}\}$$
 - [III.3.5] The scope of a declaration in a `var` statement [II.2.3] located directly within a conditional [II.2.4] or an iteration statement [II.2.5] is confined to this `var` statement.
- [III.4] The scopes [III.3.2] of any two variables sharing the same name [III.3.1] must be disjoint.
- [III.5] Any use of an identifier x [I.2.2] as an atomic expression [II.1.1] or on the left hand side of an assignment operator [II.1.7] should be in the scope [III.3] of a variable with the name x [III.3.1].
- [III.6] Every function should return a value in the following sense:
 - [III.6.1] Let us call a statement [II.2] *properly returning* if it is either a `return` statement [II.2.6], or a compound statement [II.2.2], such that the last of the statements inside it is a properly returning statement, or a conditional statement [II.2.4] with an `else` clause, such that each of the alternatives is a properly returning statement.

[III.6.2] The body of every function [II.3.2] must be a properly returning statement.

This definition of the syntax is, strictly speaking, ambiguous: the precedence of operators is undefined and the dangling else ambiguity is unresolved. However, since its sole purpose is to define the set of well-formed programs, these issues can be ignored.

The semantics of the language is left completely undefined; the set of well-formed programs gets a purely syntactical definition. It is worth note that if the semantics *were* defined in the natural way, with the domain of the sole data type set to be the set of all positive integers, then the language would be Turing-complete.

3 Boolean grammars

This section gives a quick review of *Boolean grammars*, a generalization of context-free grammars that will be used to give a formal specification of the syntax of the model programming language defined in the previous section.

A Boolean grammar [10] is a quadruple $G = (\Sigma, N, P, S)$, where Σ and N are disjoint finite nonempty sets of terminal and nonterminal symbols respectively; P is a finite set of rules of the form

$$A \rightarrow \alpha_1 \& \dots \& \alpha_m \& \neg \beta_1 \& \dots \& \neg \beta_n \quad (m + n \geq 1, \alpha_i, \beta_i \in (\Sigma \cup N)^*), \quad (1)$$

while $S \in N$ is the start symbol of the grammar. For each rule (1), the objects $A \rightarrow \alpha_i$ and $A \rightarrow \neg \beta_j$ (for all i, j) are called *conjuncts* positive and negative respectively.

A Boolean grammar is called a conjunctive grammar [7], if negation is never used, i.e., $n = 0$ for every rule (1). It degrades to a familiar context-free grammar if neither negation nor conjunction are allowed, i.e., $m = 1$ and $n = 0$ for all rules.

Multiple rules $A \rightarrow \varphi_1, \dots, A \rightarrow \varphi_k$ for a single nonterminal will be written using the standard notation

$$A \rightarrow \varphi_1 \mid \dots \mid \varphi_k$$

The semantics of Boolean grammars is defined using language equations, similarly to the well-known characterization of the context-free grammars by another simpler class language equations [4]. Let $G = (\Sigma, N, P, S)$ be a Boolean grammar. The system of language equations associated with G is a resolved system of language equations over Σ in variables N , in which the equation for each variable $A \in N$ is

$$A = \bigcup_{A \rightarrow \alpha_1 \& \dots \& \alpha_m \& \neg \beta_1 \& \dots \& \neg \beta_n \in P} \left[\bigcap_{i=1}^m \alpha_i \cap \bigcap_{j=1}^n \overline{\beta_j} \right] \quad (\text{for all } A \in N) \quad (2)$$

A system (2) can have no solutions or multiple pairwise incomparable solutions; this leads to substantial problems with the formal definition of the semantics of Boolean

grammars of the general form. These problems have been resolved using quite an elaborate technique [10]; however, in the practical cases (which exclude artistry like using a nonterminal A with a rule $A \rightarrow \neg A \& B$), a system associated with a grammar will have a unique solution (L_1, \dots, L_n) , and the language $L_G(A_i)$ generated by every i -th nonterminal A can be defined as L_i then. The language of the grammar is $L(G) = L_G(S)$.

Let us give several examples of Boolean grammars for abstract languages. The following grammar [7] employs intersection to generate $\{a^n b^n c^n \mid n \geq 0\}$, which is precisely the language Floyd [3] used to prove the non-context-freeness of the programming languages:

$$\begin{aligned} S &\rightarrow AB \& DC \\ A &\rightarrow aA \mid \varepsilon \\ B &\rightarrow bBc \mid \varepsilon \\ C &\rightarrow cC \mid \varepsilon \\ D &\rightarrow aDb \mid \varepsilon \end{aligned}$$

Note that negation was not used in this grammar, but the use of intersection is essential, as the language is not context-free.

Another non-context-free language $\{ww \mid w \in \{a, b\}^*\}$ can be denoted, this time using negation, as follows [10]:

$$\begin{aligned} S &\rightarrow \neg AB \& \neg BA \& C \\ A &\rightarrow XAX & C &\rightarrow XXC \\ A &\rightarrow a & C &\rightarrow \varepsilon \\ B &\rightarrow XBX & X &\rightarrow a \\ B &\rightarrow b & X &\rightarrow b \end{aligned}$$

A very similar language, $\{w\bar{c}w \mid w \in \{a, b\}^*\}$, can be denoted even without negation [7]:

$$\begin{aligned} S &\rightarrow C \& D \\ C &\rightarrow aCa \mid aCb \mid bCa \mid bCb \mid c \\ D &\rightarrow aA \& aD \mid bB \& bD \mid cE \\ A &\rightarrow aAa \mid aAb \mid bAa \mid bAb \mid cEa \\ B &\rightarrow aBa \mid aBb \mid bBa \mid bBb \mid cEb \\ E &\rightarrow aE \mid bE \mid \varepsilon \end{aligned}$$

The next grammar [10] defines the language $\{a^{2^n} \mid n \geq 0\}$, one more well-known non-context-free language:

$$\begin{array}{llll} S \rightarrow X \& \neg aX & X \rightarrow aX'X' & Y \rightarrow aaY'Y' & Z \rightarrow Y \\ S \rightarrow \neg X \& aX & X' \rightarrow \neg X''X'' & Y' \rightarrow \neg Y''Y'' \& T & Z \rightarrow aY \\ S \rightarrow Z \& \neg aZ & X'' \rightarrow \neg X'''X''' & Y'' \rightarrow \neg Y'''Y''' \& T & T \rightarrow aaT \\ S \rightarrow \neg Z \& aZ & X''' \rightarrow \neg X & Y''' \rightarrow \neg Y \& T & T \rightarrow \varepsilon \end{array}$$

Note that all the abstract languages considered above can be as well generated by indexed grammars [1]. However, Boolean grammars have an advantage of their own, that of being parsable in at most cubic time [10], which is as fast as the best practically usable algorithms for general context-free grammars.

Although both indexed and Boolean grammars apparently generate the same noteworthy abstract languages, their descriptive means in fact differ: one have stacks attached to nonterminals, the other have logical connectives. The descriptive machinery native to Boolean grammars shall now be used to specify the syntax of a programming language.

4 A Boolean grammar for the language

Let the alphabet of terminal symbols be $\Sigma = \{a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, _, "(", ")", \{, \}, ", ";, +, -, *, /, \%, "\&", "|", !, =, <, >\}$. This is the alphabet of the model programming language [I.1], in which the space and the newline characters are merged into a single symbol $_$, reflecting the fact that the definition of the language makes no distinction between these symbols [I.3.2].

The promised Boolean grammar over this alphabet will now be constructed. It is convenient to classify its nonterminals into the following five groups:

The start symbol $\{S\}$ corresponds to the notion of a well-defined program.

The alphabet group defines basic types of strings over the alphabet that are referred to in the main grammar. There are 11 nonterminals in this group: anyletter, anydigit, anyletterdigit, anypunctuator, anypunctuatorexceptrightpar, anypunctuatorexceptbraces, anychar, anystring, anyletterdigits, anystringwithoutbraces, safeendingstring.

The lexical group defines all types of lexemes and handles the whitespace [I], something that is typically done by a finite-automaton-based lexical analyzer. The nonterminals in this group are: WS, tVar, tIf, tElse, tWhile, tReturn, Keyword, tId, tId-aux, tNum, tNum-aux, tPlus, tMinus, tStar, tSlash, tMod, tAnd, tOr, tLessThan, tGreaterThan, tLessEqual, tGreaterEqual, tEqual, tNotEqual, BinaryOp, tUnaryMinus, tNot, UnaryOp, tLeftPar, tRightPar, tLeftBrace, tRightBrace, tAssign, tSemicolon, tComma – altogether 35 nonterminals.

The main syntax group contains nonterminals that represent the main syntactical notions of the language [III]. These are the nonterminals **Expr ExprFunctionCall, ListOfExpr, ListOfExpr1, Statement, ExprSt, CompoundSt, Statements, VarSt, CondSt, IterationSt, ReturnSt, Function, FunctionHeader, FunctionArguments, ListOfIds** and **Functions** – 17 in total.

Identifier comparison group defines a nonterminal C that generates the language

$$\{wxw \mid w \text{ is a nonempty string that consists of letters and digits;} \\ x \text{ is nonempty, it starts and ends with anything but letters or digits}\} \quad (3)$$

The following auxiliary nonterminals are used to define this language: C_{len} , $C_{iterate}$, C_{mid} , $C_a, \dots, C_z, C_0, \dots, C_9$. In total, there are $4 + 36 = 40$ nonterminals in this group.

Advanced syntax group consists of several independent components that check different syntactical constraints [III]. It contains the following 19 nonterminals: *duplicate-functions*, *duplicate-functions-here*, *all-functions-defined*, *all-functions-defined-skip*, *check-function-call*, *check-function-call-here*, *check-name*, *check-number-of-arguments*, *n-of-arg-match*, *all-variables-defined*, *all-variables-defined-skip*, *this-variable-defined*, *this-variable-defined2*, *this-variable-defined3*, *this-variable-defined3-afterskip*, *this-variable-defined4*, *these-variables-not-defined*, *has-main-function* and *returns-a-value*.

The names of nonterminals from different groups are distinguished with different fonts. Define $N = \{S\} \cup N_{alph} \cup N_{lex} \cup N_{main} \cup N_{ucw} \cup N_{adv}$. The cardinality of this set is $|N| = 1 + 11 + 35 + 17 + 40 + 19 = 123$.

Here are the simple rules for the nonterminals from the alphabet group:

$$\begin{aligned} \text{anyletter} &\rightarrow \mathbf{a} \mid \dots \mid \mathbf{z} \\ \text{anydigit} &\rightarrow \mathbf{0} \mid \dots \mid \mathbf{9} \\ \text{anyletterdigit} &\rightarrow \text{anyletter} \mid \text{anydigit} \\ \text{anypunctuator} &\rightarrow \text{“(”} \mid \text{”} \mid \{ \mid \} \mid \text{“,”} \mid \text{“;”} \mid + \mid - \mid * \mid / \mid \text{“\&”} \mid \text{“|”} \mid ! \mid \\ &= \mid < \mid > \mid \% \\ \text{anypunctuatorexceptrightpar} &\rightarrow \text{“(”} \mid \{ \mid \} \mid \text{“,”} \mid \text{“;”} \mid + \mid - \mid * \mid / \mid \text{“\&”} \mid \\ &\text{“|”} \mid ! \mid = \mid < \mid > \mid \% \\ \text{anypunctuatorexceptbraces} &\rightarrow \text{“(”} \mid \text{”} \mid \text{“,”} \mid \text{“;”} \mid + \mid - \mid * \mid / \mid \text{“\&”} \mid \\ &\text{“|”} \mid ! \mid = \mid < \mid > \mid \% \\ \text{anychar} &\rightarrow _ \mid \text{anyletter} \mid \text{anydigit} \mid \text{anypunctuator} \\ \text{anystring} &\rightarrow \text{anystring anychar} \mid \varepsilon \\ \text{anyletterdigits} &\rightarrow \text{anyletterdigits anyletterdigit} \mid \varepsilon \\ \text{anystringwithoutbraces} &\rightarrow \text{anystringwithoutbraces} _ \mid \\ &\text{anystringwithoutbraces anyletter} \mid \\ &\text{anystringwithoutbraces anydigit} \mid \\ &\text{anystringwithoutbraces anypunctuatorexceptbraces} \mid \\ &\varepsilon \\ \text{safeendingstring} &\rightarrow \text{anystring anypunctuator} \mid \text{anystring} _ \mid \varepsilon \end{aligned}$$

Rules for the lexical group do what is usually done by lexical analyzers. First, define the whitespace:

$$\text{WS} \rightarrow \text{WS } _ | \varepsilon \quad [I.3.2]$$

Now for every type of lexemes define the set of valid lexemes of this type possibly followed by whitespace.

$$\begin{aligned} \text{Keyword} &\rightarrow \text{tVar} | \text{tIf} | \text{tElse} | \text{tWhile} | \text{tReturn} & [I.2.1] \\ \text{tVar} &\rightarrow \text{v a r WS} \\ \text{tIf} &\rightarrow \text{i f WS} \\ \text{tElse} &\rightarrow \text{e l s e WS} \\ \text{tWhile} &\rightarrow \text{w h i l e WS} \\ \text{tReturn} &\rightarrow \text{r e t u r n WS} \end{aligned}$$

Let us use the extended descriptive power of Boolean grammars generously, and afford the luxury of specifying the set of identifiers precisely according to its definition:

$$\begin{aligned} \text{tId} &\rightarrow \text{tId-aux WS} \ \& \ \neg\text{Keyword} & [I.2.2] \\ \text{tId-aux} &\rightarrow \text{anyletter} | \text{tId-aux anyletter} | \text{tId-aux anydigit} \end{aligned}$$

This could have been done by simply simulating a fairly small finite automaton, though. The rest of the rules for the lexical group are dull:

$$\text{tNum} \rightarrow \text{tNum-aux WS} \quad [I.2.3]$$

$$\begin{aligned} \text{tNum-aux} &\rightarrow \text{tNum-aux anydigit} | \text{anydigit} \\ \text{BinaryOp} &\rightarrow \text{tPlus} | \text{tMinus} | \text{tStar} | \text{tSlash} | \text{tMod} | \text{tAnd} | \\ &\quad \text{tOr} | \text{tLessThan} | \text{tGreaterThan} | \text{tLessEqual} | \\ &\quad \text{tGreaterEqual} | \text{tEqual} | \text{tNotEqual} & [I.2.4] \end{aligned}$$

$$\text{tPlus} \rightarrow + \text{ WS}$$

$$\text{tMinus} \rightarrow - \text{ WS}$$

$$\text{tStar} \rightarrow * \text{ WS}$$

$$\text{tSlash} \rightarrow / \text{ WS}$$

$$\text{tMod} \rightarrow \% \text{ WS}$$

$$\text{tAnd} \rightarrow \text{“\&”} \text{ WS}$$

$$\text{tOr} \rightarrow \text{“|”} \text{ WS}$$

$$\text{tLessThan} \rightarrow < \text{ WS}$$

$$\text{tGreaterThan} \rightarrow > \text{ WS}$$

$$\text{tLessEqual} \rightarrow < = \text{ WS}$$

$$\text{tGreaterEqual} \rightarrow > = \text{ WS}$$

$$\text{tEqual} \rightarrow = = \text{ WS}$$

$$\text{tNotEqual} \rightarrow ! = \text{ WS}$$

$$\text{UnaryOp} \rightarrow \text{tUnaryMinus} | \text{tNot} \quad [I.2.5]$$

$$\text{tUnaryMinus} \rightarrow - \text{ WS}$$

$$\text{tNot} \rightarrow ! \text{ WS}$$

$$\text{tAssign} \rightarrow = \text{ WS} \quad [I.2.6]$$

$$\text{tLeftPar} \rightarrow \text{“(”} \text{ WS} \quad [I.2.7]$$

$\text{tRightPar} \rightarrow \text{"})" WS}$
 $\text{tLeftBrace} \rightarrow \{ WS$ [I.2.8]
 $\text{tRightBrace} \rightarrow \} WS$
 $\text{tSemicolon} \rightarrow \text{";" WS}$ [I.2.9]
 $\text{tComma} \rightarrow \text{"," WS}$

Next there comes the syntax group, which is composed mostly like a typical context-free grammar that partially defines the syntax of a programming language. However, note how the nonterminal **Function** invokes nonterminals from the advanced syntax group to check the variable scoping rules [III.3] and the placement of return statements [III.6].

$\text{Expr} \rightarrow \text{tId}$ [II.1.1]
 $\text{Expr} \rightarrow \text{tNum}$ [II.1.2]
 $\text{Expr} \rightarrow \text{tLeftPar Expr tRightPar}$ [II.1.3]
 $\text{Expr} \rightarrow \text{ExprFunctionCall}$ [II.1.4]
 $\text{ExprFunctionCall} \rightarrow \text{tId tLeftPar ListOfExpr tRightPar}$
 $\text{ListOfExpr} \rightarrow \text{ListOfExpr1}$
 $\text{ListOfExpr} \rightarrow \varepsilon$
 $\text{ListOfExpr1} \rightarrow \text{ListOfExpr1 tComma Expr} \mid \text{Expr}$
 $\text{Expr} \rightarrow \text{Expr BinaryOp Expr}$ [II.1.5]
 $\text{Expr} \rightarrow \text{UnaryOp Expr}$ [II.1.6]
 $\text{Expr} \rightarrow \text{tId tAssign Expr}$ [II.1.7]

$\text{Statement} \rightarrow \text{ExprSt} \mid \text{CompoundSt} \mid \text{VarSt} \mid \text{CondSt} \mid$
 $\quad \text{IterationSt} \mid \text{ReturnSt}$ [II.2]
 $\text{ExprSt} \rightarrow \text{Expr tSemicolon}$ [II.2.1]
 $\text{CompoundSt} \rightarrow \text{tLeftBrace Statements tRightBrace}$ [II.2.2]
 $\text{Statements} \rightarrow \text{Statements Statement} \mid \varepsilon$
 $\text{VarSt} \rightarrow \text{tVar} _ \text{ListOfIds tSemicolon}$ [II.2.3]
 $\text{CondSt} \rightarrow \text{tIf tLeftPar Expr tRightPar Statement} \mid$ [II.2.4]
 $\quad \text{tIf tLeftPar Expr tRightPar Statement tElse Statement}$
 $\text{IterationSt} \rightarrow \text{tWhile tLeftPar Expr tRightPar Statement}$ [II.2.5]
 $\text{ReturnSt} \rightarrow \text{tReturn Expr tSemicolon}$ [II.2.6]

$\text{FunctionHeader} \rightarrow \text{tId FunctionArguments}$ [II.3.1]
 $\text{FunctionArguments} \rightarrow \text{tLeftPar ListOfIds tRightPar} \mid$
 $\quad \text{tLeftPar tRightPar}$
 $\text{ListOfIds} \rightarrow \text{ListOfIds tComma tId} \mid \text{tId}$
 $\text{Function} \rightarrow \text{FunctionHeader CompoundSt} \&$ [II.3.2]
 $\quad \text{tId tLeftPar all-variables-defined} \&$ [III.4] [III.5]
 $\quad \text{FunctionHeader returns-a-value}$ [III.6]

$\text{Functions} \rightarrow \text{Functions Function} \mid \varepsilon$ [II.4]

The identifier comparison group defines the language (3), which is extensively used by the advanced syntax group. These rules implement the idea of the grammar for the abstract language $\{wcv \mid w \in \{a, b\}^*\}$ [7], which is reproduced in Section 3.

$$\begin{aligned}
C &\rightarrow \text{Clen} \ \& \ \text{Citerate} \mid C \ \text{WS} \\
\text{Clen} &\rightarrow \text{anyletterdigit} \ \text{Clen} \ \text{anyletterdigit} \mid \\
&\quad \text{anyletterdigit} \ \text{Cmid} \ \text{anyletterdigit} \\
C_a &\rightarrow \text{anyletterdigit} \ C_a \ \text{anyletterdigit} \mid \mathbf{a} \ \text{anyletterdigits} \ \text{Cmid} \\
&\quad \vdots \\
C_z &\rightarrow \text{anyletterdigit} \ C_z \ \text{anyletterdigit} \mid \mathbf{z} \ \text{anyletterdigits} \ \text{Cmid} \\
C_0 &\rightarrow \text{anyletterdigit} \ C_0 \ \text{anyletterdigit} \mid \mathbf{0} \ \text{anyletterdigits} \ \text{Cmid} \\
&\quad \vdots \\
C_9 &\rightarrow \text{anyletterdigit} \ C_9 \ \text{anyletterdigit} \mid \mathbf{9} \ \text{anyletterdigits} \ \text{Cmid} \\
\text{Citerate} &\rightarrow C_a \ \mathbf{a} \ \& \ \text{Citerate} \ \mathbf{a} \mid \dots \mid C_z \ \mathbf{z} \ \& \ \text{Citerate} \ \mathbf{z} \mid \\
&\quad C_0 \ \mathbf{0} \ \& \ \text{Citerate} \ \mathbf{0} \mid \dots \mid C_9 \ \mathbf{9} \ \& \ \text{Citerate} \ \mathbf{9} \mid \\
&\quad \text{anyletterdigits} \ \text{Cmid} \\
\text{Cmid} &\rightarrow \text{anypunctuator} \ \text{anystring} \ \text{anypunctuator} \mid _ \ \text{anystring} \ _ \mid \\
&\quad _ \ \text{anystring} \ \text{anypunctuator} \mid \text{anypunctuator} \ \text{anystring} \ _ \mid \\
&\quad \text{anypunctuator} \mid _
\end{aligned}$$

Finally we reach the advanced syntax group, which specifies the rules [III] that are traditionally checked by manually written program code rather than using language specification formalisms of any kind. Most of these syntactical conditions are inherently non-context-free.

Using the negation native to Boolean grammars, the rule on the uniqueness of function definition [III.1.1] can be most naturally checked by first expressing the condition that a program *has* multiply declared functions and then by putting \neg over it. The existence of duplicate function declarations is specified as follows:

$$\begin{aligned}
\text{duplicate-functions} &\rightarrow \mathbf{Functions} \ \text{duplicate-functions-here} \ \mathbf{Functions} \\
\text{duplicate-functions-here} &\rightarrow C \ \mathbf{FunctionArguments} \ \mathbf{CompoundSt}
\end{aligned}$$

A group of functions where this error occurs is “chosen” by the nonterminal *duplicate-functions*, while *duplicate-functions-here* invokes the nonterminal *C* from the identifier comparison group to ensure that the names of the first and the last functions in this group coincide. Thus a text of a program will be generated by *duplicate-functions* only if two or more functions sharing the same name are found. This is explained in Figure 1.

The requirements on the function `main` can be specified by simple context-free rules:

$$\begin{aligned}
\text{has-main-function} &\rightarrow \mathbf{Functions} \\
&\quad \mathbf{main} \ \text{tLeftPar} \ \text{tId} \ \text{tRightPar} \ \mathbf{CompoundSt} \\
&\quad \mathbf{Functions} \ [III.1.2]
\end{aligned}$$

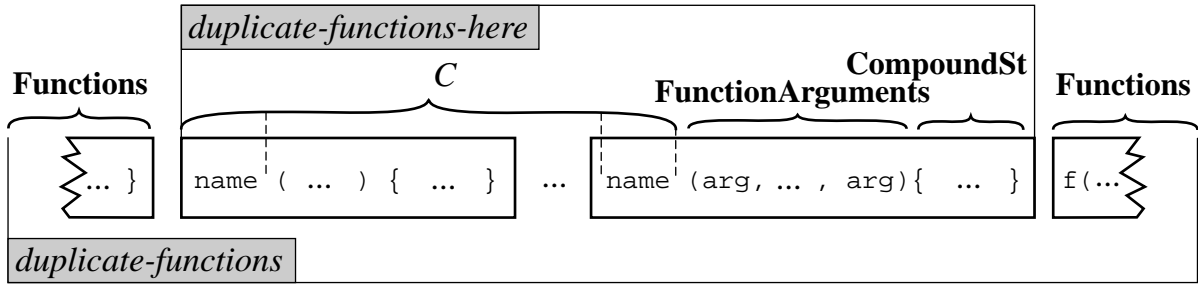


Figure 1: How *duplicate-functions* finds a pair of declarations.

Checking the correctness of all function calls [III.2] requires entering function bodies and analyzing the program lexeme by lexeme and even symbol by symbol.

$$\begin{aligned}
 \textit{all-functions-defined} &\rightarrow \textit{all-functions-defined} _ | \\
 &\quad \textit{all-functions-defined} \textit{anypunctuatorexceptrightpar} | \\
 &\quad \textit{all-functions-defined-skip} \textit{Keyword} | \\
 &\quad \textit{all-functions-defined-skip} \textit{tId} | \\
 &\quad \textit{all-functions-defined-skip} \textit{tNum} | \\
 &\quad \textit{check-function-call} \ \& \ \textit{all-functions-defined} \ \textit{tRightPar} | \\
 &\quad \varepsilon \\
 \textit{all-functions-defined-skip} &\rightarrow \textit{all-functions-defined} \ \& \ \textit{safeendingstring} \\
 \textit{check-function-call} &\rightarrow \mathbf{Functions} \ \mathbf{FunctionHeader} | \\
 &\quad \textit{anystring} \ \& \ \neg\textit{safeendingstring} \ \mathbf{ExprFunctionCall} | \\
 &\quad \mathbf{Functions} \ \textit{check-function-call-here} \\
 \textit{check-function-call-here} &\rightarrow \textit{check-name} \ \& \ \textit{check-number-of-arguments} \\
 \textit{check-name} &\rightarrow C \ \textit{tLeftPar} \ \mathbf{ListOfExpr} \ \textit{tRightPar} \quad [III.2.1] \\
 \textit{check-number-of-arguments} &\rightarrow \textit{tId} \ \textit{tLeftPar} \\
 &\quad \textit{n-of-arg-match} \ \textit{tRightPar} \quad [III.2.2] \\
 \textit{n-of-arg-match} &\rightarrow \textit{tId} \ \textit{n-of-arg-match} \ \mathbf{Expr} | \\
 &\quad \textit{tComma} \ \textit{n-of-arg-match} \ \textit{tComma} | \\
 &\quad \textit{tRightPar} \ \textit{anystring} \ \textit{tLeftPar}
 \end{aligned}$$

The nonterminal *all-functions-defined* considers all prefixes of the program and recognizes function calls [II.1.4] at the right end of the strings being considered. For each function call thus found, the nonterminal *check-function-call* is invoked, which attempts to match the call to a preceding declaration of a function. While the nonterminal *check-function-call* itself merely considers candidates for being such a declaration, another nonterminal *check-function-call-here* does the actual matching (see Figure 2), using the nonterminal *check-name* to ensure that the function defined at the beginning in the current substring is the same function as the one called at the end of the current substring, and at the same time the nonterminal *check-number-of-arguments* checks that

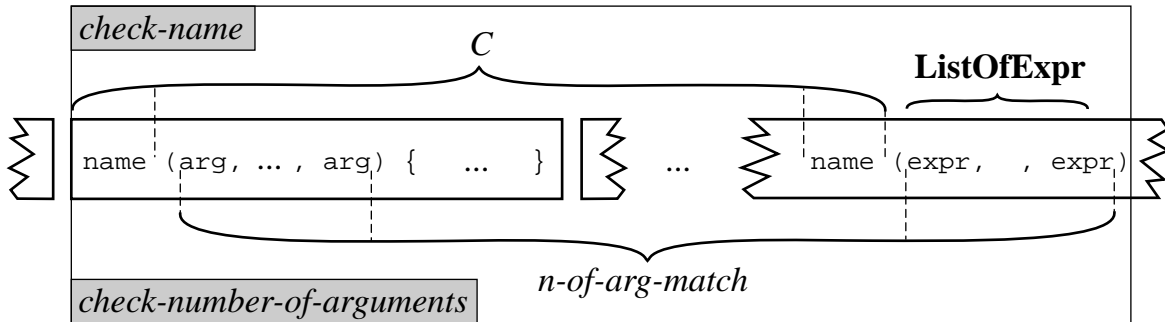


Figure 2: How *check-function-call-here* matches a declaration [II.3.1] to a use [II.1.4].

the number of arguments in the call corresponds to the number of formal arguments in the declaration.

Now consider the rules for the nonterminal *all-variables-defined*, which has already been used in the rules for the syntax group.

```

all-variables-defined → all-variables-defined _ |
  all-variables-defined anypunctuator & ¬safeendingstring VarSt |
  all-variables-defined-skip tId tLeftPar |
  all-variables-defined-skip tNum |
  all-variables-defined-skip Keyword |
  these-variables-not-defined tSemicolon &
    all-variables-defined-skip VarSt |
  this-variable-defined & all-variables-defined-skip tId |
  these-variables-not-defined tRightPar
all-variables-defined-skip → all-variables-defined & safeendingstring
this-variable-defined → C | tId tComma this-variable-defined |
  tId tRightPar tLeftBrace this-variable-defined2 |
this-variable-defined2 → Statement this-variable-defined2 |
  tLeftBrace this-variable-defined2 |
  tIf tLeftPar Expr tRightPar
    tLeftBrace this-variable-defined2 |
  tIf tLeftPar Expr tRightPar Statement tElse
    tLeftBrace this-variable-defined2 |
  tWhile tLeftPar Expr tRightPar
    tLeftBrace this-variable-defined2 |
  tVar _ this-variable-defined3 |
  tIf tLeftPar Expr tRightPar tVar _ this-variable-defined4 |
  tIf tLeftPar Expr tRightPar Statement tElse
    tVar _ this-variable-defined4 |
  tWhile tLeftPar Expr tRightPar tVar _ this-variable-defined4

```

$$\begin{aligned}
\textit{this-variable-defined}_3 &\rightarrow \text{tId } \text{tComma } \textit{this-variable-defined}_3 \mid \\
&\quad C \ \& \ \text{anystringwithoutbraces } \textit{this-variable-defined}_3\text{-afterskip} \\
\textit{this-variable-defined}_3\text{-afterskip} &\rightarrow \text{anystringwithoutbraces} \\
&\quad \mathbf{Statement } \textit{this-variable-defined}_3\text{-afterskip} \mid \\
&\quad \text{tLeftBrace } \textit{this-variable-defined}_3\text{-afterskip} \mid \\
\textit{this-variable-defined}_4 &\rightarrow \text{tId } \text{tComma } \textit{this-variable-defined}_4 \mid \\
&\quad C \ \& \ \mathbf{ListOfIds};
\end{aligned}$$

The nonterminal *all-variables-defined* considers all prefixes of the body of a function and recognizes declarations statements [II.2.3] and references to variables [III.1.1] at the right ends of these prefixes:

- Whenever an atomic expression [II.1.1] is found at the right, the nonterminal *this-variable-defined* is used to match it to a declaration of a variable [III.3] at the left. First it tries to find such a declaration among the formal arguments of the function [III.3.3]. If failed, the nonterminal *this-variable-defined*₂ starts looking for a declaration in a **var** statement, observing the scoping rules. **var** statements located inside compound statements [III.3.4] are handled by the nonterminal *this-variable-defined*₃, while *this-variable-defined*₄ takes care of declarations directly inside **if** and **while** statements [III.3.5].
- When a declaration statement [II.2.3] appears at the right, the nonterminal *these-variables-not-defined* is used to ensure that none of the declared variables is in the scope of a variables with the same name.

Consider the rules for *these-variables-not-defined*. Here negation comes particularly handy, since the condition that an identifier is in the scope of a variable with the same name has already been expressed by the nonterminal *this-variable-defined*, and it suffices to invert it.

$$\begin{aligned}
\textit{these-variables-not-defined} &\rightarrow \textit{these-variables-not-defined} \text{tComma } \text{tId} \ \& \\
&\quad \neg \textit{this-variable-defined} \mid \\
&\quad \text{safeendingstring } \text{tVar } _ \text{tId} \ \& \ \neg \textit{this-variable-defined} \mid \\
&\quad \text{tId} \mid \varepsilon
\end{aligned}$$

The condition that a function returns a value can be checked by simple context-free rules as follows:

$$\begin{aligned}
\textit{returns-a-value} &\rightarrow \text{tLeftBrace } \mathbf{Statements} \textit{returns-a-value} \text{tRightBrace} \mid \\
&\quad \text{tIf } \text{tLeftPar } \mathbf{Expr} \text{tRightPar } \textit{returns-a-value} \\
&\quad \quad \text{tElse } \textit{returns-a-value} \mid \\
&\quad \mathbf{ReturnSt}
\end{aligned}$$

Finally, a single rule for the start symbol specifies what a well-formed program is:

S → WS Functions &	[II.4]
WS <i>all-functions-defined</i> &	[III.2]
WS <i>has-main-function</i> &	[III.1.2]
¬WS <i>duplicate-functions</i>	[III.1.1]

This completes the grammar, which consists of 123 nonterminals and 368 rules, and specifies the language as defined in Section 2. Since it is known that every language generated by a Boolean grammar can be parsed in time $O(n^3)$, this grammar directly yields a polynomial-time correctness checker for the model programming language.

The given grammar has in fact been extensively tested using a parser generator [9], which implements a yet unpublished extension of the Generalized LR parsing algorithm [13, 8] for Boolean grammars [12]. The observed parsing time was $O(n^2)$. The details of the tests and the actual performance results are given in Appendixes A through C.

5 Conclusion

A specification of a programming language by a formal grammar that belongs a computationally feasible class of grammars was obtained, apparently for the first time ever.

The given Boolean grammar for the model language was admittedly composed under an influence of the standard compiler construction methods that separate lexical analysis, context-free syntax backbone and non-context-free syntactical restrictions [2]. Perhaps a much better grammar for the same model language could be constructed using some other principles that would make better use of the expanded descriptive means of Boolean grammars. Realizing such principles and coming up with improved grammar composition techniques might be of interest.

Another obvious topic for further study would be to investigate the formal properties of Boolean grammars [10] and to develop parsing algorithms for them. Some algorithms are already being researched [11, 12], and their prototype implementations are available [9]. It will require more work to make these algorithms practical.

One more direction for future work is to use the experience of this study and come up with a new grammar formalism with better properties, which would still allow to specify simple programming languages. That might be an efficiently parsable and still sufficiently expressive subclass of Boolean grammars, or that might be an entirely new formalism.

In any case, it was shown that formal grammars are better suited for specifying the syntax of programming languages than it has been believed for the last three decades, and thus deserve renewed attention.

References

- [1] A. V. Aho, “Indexed grammars”, *Journal of the ACM*, 15:4 (1968), 647–671.
- [2] A. V. Aho, R. Sethi, J. D. Ullman, *Compilers: principles, techniques and tools*, Addison-Wesley, Reading, Mass., 1986.
- [3] R. W. Floyd, “On the non-existence of a phrase structure grammar for ALGOL 60”, *Communications of the ACM*, 5 (1962), 483–484.
- [4] S. Ginsburg, H. G. Rice, “Two families of languages related to ALGOL”, *Journal of the ACM*, 9 (1962), 350–371.
- [5] P. Naur (Ed.), J. W. Backus, J. H. Wegstein, A. van Wijngaarden, M. Woodger, F. L. Bauer, J. Green, C. Katz, J. McCarthy, A. J. Perlis, H. Rutishauser, K. Samelson, B. Vauquois, “Report on the algorithmic language ALGOL 60”, *Communications of the ACM*, 3:5 (1960), 299–314.
- [6] P. Naur (Ed.), J. W. Backus, J. H. Wegstein, A. van Wijngaarden, M. Woodger, F. L. Bauer, J. Green, C. Katz, J. McCarthy, A. J. Perlis, H. Rutishauser, K. Samelson, B. Vauquois, “Revised report on the algorithmic language ALGOL 60”, *Communications of the ACM*, 6:1 (1963), 1–17.
- [7] A. Okhotin, “Conjunctive grammars”, *Journal of Automata, Languages and Combinatorics*, 6:4 (2001), 519–535.
- [8] A. Okhotin, “LR parsing for conjunctive grammars”, *Grammars*, 5:2 (2002), 81–124.
- [9] A. Okhotin, “Whale Calf, a parser generator for conjunctive grammars”, *Implementation and Application of Automata* (Proceedings of CIAA 2002, Tours, France, July 3–5, 2002), LNCS 2608, 213–220. The software is available at

<http://www.cs.queensu.ca/home/okhotin/whalecalf/>
- [10] A. Okhotin, “Boolean grammars”, *Information and Computation*, to appear.
 - Preliminary version in: *Developments in Language Theory* (Proceedings of DLT 2003, Szeged, Hungary, July 7–11, 2003), LNCS 2710, 398–410.
- [11] A. Okhotin, “An extension of recursive descent parsing for Boolean grammars”, Technical Report 2004–475, School of Computing, Queen’s University, Kingston, Ontario, Canada.
- [12] A. Okhotin, “LR parsing for Boolean grammars”, in progress.

- [13] M. Tomita, “An efficient augmented context-free parsing algorithm”, *Computational Linguistics*, 13:1 (1987), 31–46.
- [14] A. van Wijngaarden, “Orthogonal design and description of a formal language”, Technical Report MR 76, Mathematisch Centrum, Amsterdam, 1965.
- [15] A. van Wijngaarden, B. J. Mailloux, J. E. L. Peck, C. H. A. Koster, M. Sintzoff, C. H. Lindsey, L. G. L. T. Meertens, R. G. Fisker, “Revised report on the algorithmic language ALGOL 68”, *Acta Informatica*, 5 (1975), 1–236.

Appendix.

A A parser

Using the Whale Calf parser generator [9], a parser for the model programming language has been automatically created from the Boolean grammar given in Section 4.

The following main module was used, which simply loads the input file (translating the character set), invokes the generated parser, and afterwards prints the result of recognition.

```
mlc.cpp #include <fstream>
#include <vector>
#include "ml_parser.h"
using namespace std;

void preprocessor(istream &is, vector<int> &v);

int main(int argc, char *argv[])
{
    ifstream is(argv[1]);
    if(!is) { cerr << "Unable to open "
               << argv[1] << "\n";
              return 1; }
    vector<int> v;
    preprocessor(is, v);
    cout << "The input is " << v.size()
          << " symbols long.\n";

    WhaleCalfLRParser whale_calf(Ml_parser::grammar);
    cout << whale_calf.copyright_notice() << "\n";
    whale_calf.read(v);
    bool result=whale_calf.recognize();
    if(result)
        cout << "This is a well-formed program.\n";
    else
        cout << "Syntax errors found.\n";
    return result;
}

void preprocessor(istream &is, vector<int> &v)
{
    for(int i=0; i++)
    {
        char c=is.get();
        if(is.eof()) return;

        int x=0;
        if(c==' ' || c=='\t' || c=='\n')
            x=Ml_parser::Terminal_space;
        else if(c>='a' && c<='z')
            x=Ml_parser::Terminal_a+int(c-'a');
```

```
        else if(c>='0' && c<='9')
            x=Ml_parser::Terminal_0+int(c-'0');
        else if(c=='(')
            x=Ml_parser::Terminal_leftpar;
        else if(c==')')
            x=Ml_parser::Terminal_rightpar;
        else if(c=='{')
            x=Ml_parser::Terminal_leftbrace;
        else if(c=='}')
            x=Ml_parser::Terminal_rightbrace;
        else if(c==',')
            x=Ml_parser::Terminal_comma;
        else if(c==';')
            x=Ml_parser::Terminal_semicolon;
        else if(c=='+')
            x=Ml_parser::Terminal_plus;
        else if(c=='-')
            x=Ml_parser::Terminal_minus;
        else if(c=='*')
            x=Ml_parser::Terminal_star;
        else if(c=='/')
            x=Ml_parser::Terminal_slash;
        else if(c=='%')
            x=Ml_parser::Terminal_percent;
        else if(c=='&')
            x=Ml_parser::Terminal_and;
        else if(c=='|')
            x=Ml_parser::Terminal_or;
        else if(c=='!')
            x=Ml_parser::Terminal_excl;
        else if(c=='=')
            x=Ml_parser::Terminal_eq;
        else if(c=='<')
            x=Ml_parser::Terminal_lt;
        else if(c=='>')
            x=Ml_parser::Terminal_gt;
        else
        {
            cerr << "Invalid symbol at offset "
                  << i << ".\n";
            throw false;
        }

        v.push_back(x);
    }
}
```

The following grammar description file for Whale Calf was used to specify the parser. The grammar is identical to the one given in Section 4.

```
algorithm=LR;
terminal_space;
```

ml.whc

```

terminal _a, _b, _c, _d, _e, _f, _g, _h, _i, _j,
    _k, _l, _m, _n, _o, _p, _q, _r, _s, _t,
    _u, _v, _w, _x, _y, _z;
terminal _0, _1, _2, _3, _4, _5, _6, _7, _8, _9;
terminal _leftpar, _rightpar, _leftbrace,
    _rightbrace, _comma, _semicolon;
terminal _plus, _minus, _star, _slash, _percent,
    _and, _or, _excl, _lt, _gt;

S -> WS Functions &
    WS all_functions_defined &
    WS has_main_function &
    ~WS duplicate_functions;

WS -> WS _space | e;
anyletter -> _a | _b | _c | _d | _e | _f | _g | _h |
    _i | _j | _k | _l | _m | _n | _o | _p | _q |
    _r | _s | _t | _u | _v | _w | _x | _y | _z;
anydigit -> _0 | _1 | _2 | _3 | _4 | _5 | _6 | _7 |
    _8 | _9;
anyletterdigit -> anyletter | anydigit;
anypunctuator -> _leftpar | _rightpar | _leftbrace |
    _rightbrace | _comma | _semicolon | _plus |
    _minus | _star | _slash | _and | _or | _excl |
    _eq | _lt | _gt | _percent;
anypunctuatorexceptrightpar -> _leftpar |
    _leftbrace | _rightbrace | _comma | _semicolon |
    _plus | _minus | _star | _slash | _and | _or |
    _excl | _eq | _lt | _gt | _percent;
anypunctuatorexceptbraces -> _leftpar | _rightpar |
    _comma | _semicolon | _plus | _minus | _star |
    _slash | _and | _or | _excl | _eq | _lt | _gt |
    _percent;
anychar -> _space | anyletter | anydigit |
    anypunctuator;
anystring -> anystring anychar | e;
anyletterdigits -> anyletterdigits anyletterdigit |
    e;
anystringwithoutbraces ->
    anystringwithoutbraces _space |
    anystringwithoutbraces anyletter |
    anystringwithoutbraces anydigit |
    anystringwithoutbraces
        anypunctuatorexceptbraces |
    e;

tVar -> _v _a _r WS;
tIf -> _i _f WS;
tElse -> _e _l _s _e WS;
tWhile -> _w _h _i _l _e WS;
tReturn -> _r _e _t _u _r _n WS;
Keyword -> tVar | tIf | tElse | tWhile | tReturn;

tId -> tId_aux WS & ~Keyword;
tId_aux -> anyletter | tId_aux anyletter |
    tId_aux anydigit;
tNum -> tNum_aux WS;
tNum_aux -> tNum_aux anydigit | anydigit;

ListOfIds -> ListOfIds tComma tId | tId;

tPlus -> _plus WS;

tMinus -> _minus WS;
tStar -> _star WS;
tSlash -> _slash WS;
tMod -> _percent WS;
tAnd -> _and WS;
tOr -> _or WS;
tLessThan -> _lt WS;
tGreaterThan -> _gt WS;
tLessEqual -> _lt _eq WS;
tGreaterEqual -> _gt _eq WS;
tEqual -> _eq _eq WS;
tNotEqual -> _excl _eq WS;
BinaryOp -> tPlus | tMinus | tStar | tSlash | tMod |
    tAnd | tOr | tLessThan | tGreaterThan |
    tLessEqual | tGreaterEqual | tEqual | tNotEqual;

tUnaryMinus -> _minus WS;
tNot -> _excl WS;
UnaryOp -> tUnaryMinus | tNot;

tLeftPar -> _leftpar WS;
tRightPar -> _rightpar WS;
tLeftBrace -> _leftbrace WS;
tRightBrace -> _rightbrace WS;
tAssign -> _eq WS;
tSemicolon -> _semicolon WS;
tComma -> _comma WS;

Expr -> tId;
Expr -> tNum;
Expr -> tLeftPar Expr tRightPar;
Expr -> ExprFunctionCall;
Expr -> Expr BinaryOp Expr;
Expr -> UnaryOp Expr;
Expr -> tId tAssign Expr;

ExprFunctionCall ->
    tId tLeftPar ListOfExpr tRightPar;
ListOfExpr -> ListOfExpr1;
ListOfExpr -> e;
ListOfExpr1 -> ListOfExpr1 tComma Expr | Expr;

Statement -> ExprSt | CompoundSt | VarSt | CondSt |
    IterationSt | ReturnSt;
ExprSt -> Expr tSemicolon;
CompoundSt -> tLeftBrace Statements tRightBrace;
Statements -> Statements Statement | e;
VarSt -> tVar _space ListOfIds tSemicolon;
CondSt -> tIf tLeftPar Expr tRightPar Statement |
    tIf tLeftPar Expr tRightPar Statement
        tElse Statement;
IterationSt -> tWhile tLeftPar Expr tRightPar
    Statement;
ReturnSt -> tReturn Expr tSemicolon;

Function -> FunctionHeader CompoundSt &
    tId tLeftPar all_variables_defined &
    FunctionHeader returns_a_value;
FunctionHeader -> tId FunctionArguments;
FunctionArguments -> tLeftPar ListOfIds tRightPar |
    tLeftPar tRightPar;

```

```

Functions -> Functions Function | e;

safeendingstring -> anystring anypunctuator |
  anystring _space | e;

C -> Clen & Citerate | C WS;
Clen -> anyletterdigit Clen anyletterdigit |
  anyletterdigit Cmid anyletterdigit;
C_a -> anyletterdigit C_a anyletterdigit |
  _a anyletterdigits Cmid;
C_b -> anyletterdigit C_b anyletterdigit |
  _b anyletterdigits Cmid;
C_c -> anyletterdigit C_c anyletterdigit |
  _c anyletterdigits Cmid;
C_d -> anyletterdigit C_d anyletterdigit |
  _d anyletterdigits Cmid;
C_e -> anyletterdigit C_e anyletterdigit |
  _e anyletterdigits Cmid;
C_f -> anyletterdigit C_f anyletterdigit |
  _f anyletterdigits Cmid;
C_g -> anyletterdigit C_g anyletterdigit |
  _g anyletterdigits Cmid;
C_h -> anyletterdigit C_h anyletterdigit |
  _h anyletterdigits Cmid;
C_i -> anyletterdigit C_i anyletterdigit |
  _i anyletterdigits Cmid;
C_j -> anyletterdigit C_j anyletterdigit |
  _j anyletterdigits Cmid;
C_k -> anyletterdigit C_k anyletterdigit |
  _k anyletterdigits Cmid;
C_l -> anyletterdigit C_l anyletterdigit |
  _l anyletterdigits Cmid;
C_m -> anyletterdigit C_m anyletterdigit |
  _m anyletterdigits Cmid;
C_n -> anyletterdigit C_n anyletterdigit |
  _n anyletterdigits Cmid;
C_o -> anyletterdigit C_o anyletterdigit |
  _o anyletterdigits Cmid;
C_p -> anyletterdigit C_p anyletterdigit |
  _p anyletterdigits Cmid;
C_q -> anyletterdigit C_q anyletterdigit |
  _q anyletterdigits Cmid;
C_r -> anyletterdigit C_r anyletterdigit |
  _r anyletterdigits Cmid;
C_s -> anyletterdigit C_s anyletterdigit |
  _s anyletterdigits Cmid;
C_t -> anyletterdigit C_t anyletterdigit |
  _t anyletterdigits Cmid;
C_u -> anyletterdigit C_u anyletterdigit |
  _u anyletterdigits Cmid;
C_v -> anyletterdigit C_v anyletterdigit |
  _v anyletterdigits Cmid;
C_w -> anyletterdigit C_w anyletterdigit |
  _w anyletterdigits Cmid;
C_x -> anyletterdigit C_x anyletterdigit |
  _x anyletterdigits Cmid;
C_y -> anyletterdigit C_y anyletterdigit |
  _y anyletterdigits Cmid;
C_z -> anyletterdigit C_z anyletterdigit |
  _z anyletterdigits Cmid;
C_0 -> anyletterdigit C_0 anyletterdigit |
  _0 anyletterdigits Cmid;

C_1 -> anyletterdigit C_1 anyletterdigit |
  _1 anyletterdigits Cmid;
C_2 -> anyletterdigit C_2 anyletterdigit |
  _2 anyletterdigits Cmid;
C_3 -> anyletterdigit C_3 anyletterdigit |
  _3 anyletterdigits Cmid;
C_4 -> anyletterdigit C_4 anyletterdigit |
  _4 anyletterdigits Cmid;
C_5 -> anyletterdigit C_5 anyletterdigit |
  _5 anyletterdigits Cmid;
C_6 -> anyletterdigit C_6 anyletterdigit |
  _6 anyletterdigits Cmid;
C_7 -> anyletterdigit C_7 anyletterdigit |
  _7 anyletterdigits Cmid;
C_8 -> anyletterdigit C_8 anyletterdigit |
  _8 anyletterdigits Cmid;
C_9 -> anyletterdigit C_9 anyletterdigit |
  _9 anyletterdigits Cmid;

Citerate ->
  C_a_a & Citerate_a | C_b_b & Citerate_b |
  C_c_c & Citerate_c | C_d_d & Citerate_d |
  C_e_e & Citerate_e | C_f_f & Citerate_f |
  C_g_g & Citerate_g | C_h_h & Citerate_h |
  C_i_i & Citerate_i | C_j_j & Citerate_j |
  C_k_k & Citerate_k | C_l_l & Citerate_l |
  C_m_m & Citerate_m | C_n_n & Citerate_n |
  C_o_o & Citerate_o | C_p_p & Citerate_p |
  C_q_q & Citerate_q | C_r_r & Citerate_r |
  C_s_s & Citerate_s | C_t_t & Citerate_t |
  C_u_u & Citerate_u | C_v_v & Citerate_v |
  C_w_w & Citerate_w | C_x_x & Citerate_x |
  C_y_y & Citerate_y | C_z_z & Citerate_z;

Citerate ->
  C_0_0 & Citerate_0 | C_1_1 & Citerate_1 |
  C_2_2 & Citerate_2 | C_3_3 & Citerate_3 |
  C_4_4 & Citerate_4 | C_5_5 & Citerate_5 |
  C_6_6 & Citerate_6 | C_7_7 & Citerate_7 |
  C_8_8 & Citerate_8 | C_9_9 & Citerate_9;

Citerate -> anyletterdigits Cmid;
Cmid -> anypunctuator anystring anypunctuator |
  _space anystring anypunctuator |
  anypunctuator anystring _space |
  _space anystring _space | anypunctuator |
  _space;

duplicate_functions ->
  Functions duplicate_functions_here Functions;
duplicate_functions_here ->
  C FunctionArguments CompoundSt;

all_functions_defined -> all_functions_defined
  anypunctuatorexcepttrightpar |
  all_functions_defined _space |
  all_functions_defined_skip Keyword |
  all_functions_defined_skip tId |
  all_functions_defined_skip tNum |
  check_function_call &
  all_functions_defined tRightPar |
  e;
all_functions_defined_skip ->
  all_functions_defined & safeendingstring;

```

```

check_function_call -> Functions FunctionHeader |
  anystring & ~safeendingstring ExprFunctionCall |
  Functions check_function_call_here;

check_function_call_here ->
  check_name & check_number_of_arguments;
check_name -> C tLeftPar ListOfExpr tRightPar;
check_number_of_arguments ->
  tId tLeftPar n_of_arg_match tRightPar;
n_of_arg_match -> tId n_of_arg_match Expr |
  tComma n_of_arg_match tComma |
  tRightPar anystring tLeftPar;

all_variables_defined ->
  all_variables_defined _space |
  all_variables_defined anypunctuator &
  ~safeendingstring VarSt |
  all_variables_defined_skip tId tLeftPar |
  all_variables_defined_skip tNum |
  all_variables_defined_skip Keyword |
  these_variables_not_defined tSemicolon &
  all_variables_defined_skip VarSt |
  this_variable_defined &
  all_variables_defined_skip tId |
  these_variables_not_defined tRightPar;
all_variables_defined_skip ->
  all_variables_defined & safeendingstring;

this_variable_defined -> C |
  tId tComma this_variable_defined |
  tId tRightPar tLeftBrace this_variable_defined2;
this_variable_defined2 ->
  Statement this_variable_defined2 |
  tLeftBrace this_variable_defined2 |
  tIf tLeftPar Expr tRightPar tLeftBrace
  this_variable_defined2 |
  tIf tLeftPar Expr tRightPar Statement
  tElse tLeftBrace this_variable_defined2 |
  tWhile tLeftPar Expr tRightPar tLeftBrace
  this_variable_defined2 |
  tVar _space this_variable_defined3 |
  tIf tLeftPar Expr tRightPar
  tVar _space this_variable_defined4 |
  tIf tLeftPar Expr tRightPar Statement
  tElse tVar _space this_variable_defined4 |
  tWhile tLeftPar Expr tRightPar
  tVar _space this_variable_defined4;
this_variable_defined3 ->
  tId tComma this_variable_defined3 |
  C & anystringwithoutbraces
  this_variable_defined3_afterskip;
this_variable_defined3_afterskip ->
  Statement this_variable_defined3_afterskip |
  tLeftBrace this_variable_defined3_afterskip |
  anystringwithoutbraces;
this_variable_defined4 ->
  tId tComma this_variable_defined4 |
  C & ListOfIds;

these_variables_not_defined ->
  these_variables_not_defined tComma tId &
  ~this_variable_defined |
  safeendingstring tVar _space tId &
  ~this_variable_defined |
  tId | e;

has_main_function -> Functions _m_a_i_n
  tLeftPar tId tRightPar CompoundSt Functions;
returns_a_value ->
  tLeftBrace Statements
  returns_a_value tRightBrace |
  tIf tLeftPar Expr tRightPar returns_a_value
  tElse returns_a_value |
  ReturnSt;

```

The generated parser uses the Generalized LR algorithm [12]. The parsing table contains 743 states.

B Test suite

This appendix contains an extensive collection of programs in the model language, which attempts to cover all aspects of its syntax. These are the tests the author used to convince himself that the grammar is correct, and they are given here in hope to convince the reader of the same.

Each example has a file name, which ends with **-yes** if the program is well-formed, with **-no** otherwise. This name is given in a margin note that accompanies every listing. Every example was tested with the parser, and its computation time¹ and the actual result of the computation (acceptance or rejection) are included in the corresponding margin note.

The first group of tests covers the lexical component of the language [I]. Following is a well-formed program that makes use of all the symbols from the alphabet [I.1], contains at least one lexeme [I.2] of each type, and plenty of whitespace [I.3]:

```

main(arg)
{

```

I-a0-yes
3.17 sec., Acc

¹A Pentium 4/1.5 GHz system running Windows 2000 was used. The parser was built with Borland C++ 5.5.1 command-line compiler using the options `-O2 -6`.

```

var x, x1az, x1;
var qwertyuiop, asdfghjkl, z1x23c445vbn7m890;

x=x1az*100/(1+x1)-x%arg;
if(x1>x1az & !(x<x1az | x>x1az | x==x1az))
    while(2+2!=4)
        x=-x;
else
    x=0;

qwertyuiop=asdfghjkl=z1x23c445vbn7m890=x;
return x;
}

```

Let us now remove all nonessential whitespace from this program:

I-a1-yes
2.51 sec., Acc

```

main(arg){var x,x1az,x1;var qwertyuiop,asdfghjkl,z1x23c445vbn7m890;x=x1az*100/(1+x1)-x%arg;if(x1>x1az&!(x<x1az|x>x1az|x==x1az))while(2+2!=4)x=-x;else x=0;qwertyuiop=asdfghjkl=z1x23c445vbn7m890=x;return x;}

```

Only four essential spaces remain. The removal of any of them renders the program ill-formed:

I-a2-no
1.11 sec., Rej

```

main(arg){varx,x1az,x1;var qwertyuiop,asdfghjkl,z1x23c445vbn7m890;x=x1az*100/(1+x1)-x%arg;if(x1>x1az&!(x<x1az|x>x1az|x==x1az))while(2+2!=4)x=-x;else x=0;qwertyuiop=asdfghjkl=z1x23c445vbn7m890=x;return x;}

```

The next shorter example is a well-formed program:

I-b0-yes
0.43 sec., Acc

```

main(arg)
{
    var x, elsex;
    if(arg!=1) 1; else x=1;
    return 1;
}

```

Now let us remove the space between `var` and `x`, changing this pair of a keyword [I.2.1] and of an identifier [I.2.2] to a single identifier. The resulting program thus contains an incorrect statement formed by an identifier `varx`, a comma [I.2.9], an identifier and a semicolon:

I-b1-no
0.31 sec., Rej

```

main(arg)
{
    varx, elsex;
    if(arg!=1) 1; else x=1;
    return 1;
}

```

The removal of a space between the keyword `return` [I.2.1] and the number 1 [I.2.3] turns these two lexemes into a single identifier `return1` [I.2.2], which forms an atomic expression [II.1.1], valid in itself. However, as there is no variable with this name, this violates one of the syntactical rules [III.5], and also the function no longer returns a value [III.6]:

```

main(arg)
{
    var x, elsex;
    if(arg!=1) 1; else x=1;
    return1;
}

```

I-b2-no
0.42 sec., Rej

Let us try to spoil [I-b0-yes] in the same way by removing the space between the keyword `else` [I.2.1] and the identifier `else`, turning them into a single identifier `elsex`. Surprisingly, this time the program remains well-formed, as the `if` statement is still correct after losing its `else` clause [II.2.4], while `elsex=1;` is a valid expression-statement [II.2.1], because the variable `elsex` happens to be declared [III.3.1] before.

```

main(arg)
{
    var x, elsex;
    if(arg!=1) 1; elsex=1;
    return 1;
}

```

I-b3-yes
0.42 sec., Acc

Splitting lexemes in [I-b0-yes] by whitespace characters makes the program ill-formed:

```

main(arg)
{
    var x, elsex;
    if(arg! =1) 1; else x=1;
    return 1;
}

```

I-b4-no
0.42 sec., Rej

```

main(arg)
{
    var x, elsex;
    if(arg!=1) 1; else x=1;
    ret urn 1;
}

```

I-b5-no
0.43 sec., Rej

Keywords [I.2.1] do not work as identifiers [I.2.2]:

I-b6-no
0.42 sec., Rej

```
main(arg)
{
  var x, while;
  if(arg!=1) 1; else x=1;
  return 1;
}
```

Identifiers cannot start from digits:

I-b7-no
0.42 sec., Rej

```
main(arg)
{
  var x, 7elsex;
  if(arg!=1) 1; else x=1;
  return 1;
}
```

The next group of tests covers all types of syntactical constructs [II] defined in the language. Let us start from the syntax of the expressions [II.1]:

II.1-a0-yes
0.57 sec., Acc

```
f(x, y)
{
  return x+y;
}

main(arg)
{
  var x;
  x=f(arg+1, x=1)*-x/(x+arg);
  return x;
}
```

Its corrupted version with one parenthesis removed [II.1.3] is ill-formed:

II.1-a1-no
0.55 sec., Rej

```
f(x, y)
{
  return x+y;
}

main(arg)
{
  var x;
  x=f(arg+1, x=1)*-x/x+arg);
  return x;
}
```

Another corrupted version of [II.1-a0-yes] tries to use the symbol * as a prefix unary operator [II.1.6], which it isn't [I.2.5].

II.1-a2-no
0.56 sec., Rej

```
f(x, y)
{
  return x+y;
}

main(arg)
{
  var x;
  x=f(arg+1, x=1)*-x/(x+arg);
  return x;
}
```

```
{
  var x;
  x=f(arg+1, x=1)-*x/(x+arg);
  return x;
}
```

The next corrupted expression attempts a postfix use of the minus:

```
f(x, y)
{
  return x+y-;
}

main(arg)
{
  var x;
  x=f(arg+1, x=1)*-x/(x+arg);
  return x;
}
```

II.1-a3-no
0.43 sec., Rej

The following modification violates the format of function calls [II.1.4]:

```
f(x, y)
{
  return x+y;
}

main(arg)
{
  var x;
  x=f(arg+1 x=1)*-x/(x+arg);
  return x;
}
```

II.1-a4-no
0.56 sec., Rej

Assignment expression syntax [II.1.7] is not observed here:

```
f(x, y)
{
  return x+y;
}

main(arg)
{
  var x;
  x=f(arg+1, x*2=1)*-x/(x+arg);
  return x;
}
```

II.1-a5-no
0.57 sec., Rej

Now let us turn to the syntax of statements [II.2]. This well-formed program contains all types of statements defined in the language:

```
main(arg)
{
  var x;
```

II.2-a0-yes
0.73 sec., Acc


```

    if(x>x)
        while(x==x)
            x=x;
    else
    {
        if(x==arg)
            return x;
        { }
        x;
    }
    return x;
}

```

Consider a short correct program, and remove a semicolon that is a part of an expression-statement [II.2.1]:

II.2-b0-yes
0.20 sec., Acc

```

main(arg)
{
    arg=arg;
    return arg;
}

```

II.2-b1-no
0.18 sec., Rej

```

main(arg)
{
    arg=arg
    return arg;
}

```

Use a semicolon outside of a statement (there are no C-like empty statements in the language):

II.2-b2-no
0.15 sec., Rej

```

main(arg)
{
    ;
    return arg;
}

```

Use a compound statement [II.2.2], and then forget one of its delimiters:

II.2-b3-yes
0.21 sec., Acc

```

main(arg)
{
    { arg=arg; }
    return arg;
}

```

II.2-b4-no
0.20 sec., Rej

```

main(arg)
{
    { arg=arg;
    return arg;
}

```

Add a declaration [II.2.3], and then violate its syntax:

```

main(arg)
{
    var x;
    arg=arg;
    return arg;
}

```

II.2-b5-yes
0.28 sec., Acc

```

main(arg)
{
    var ;
    arg=arg;
    return arg;
}

```

II.2-b6-no
0.25 sec., Rej

Consider another small well-formed program:

```

main(arg)
{
    if(1)
        arg=arg;

    return 1;
}

```

II.2-c0-yes
0.25 sec., Acc

Supply an else clause to the conditional statement [II.2.4]:

```

main(arg)
{
    if(1)
        arg=arg;
    else
        arg=-arg;

    return 1;
}

```

II.2-c1-yes
0.36 sec., Acc

Use two else clauses, which is not allowed [II.2.4]:

```

main(arg)
{
    if(1)
        arg=arg;
    else
        arg=-arg;
    else
        arg=arg+arg;

    return 1;
}

```

II.2-c2-no
0.45 sec., Rej

Use an else clause outside a conditional statement [II.2.4]:

```

II.2-c3-no main(arg)
0.23 sec., Rej {
    else
        arg=arg;

    return 1;
}

```

Omit the parentheses:

```

II.2-c4-no main(arg)
0.21 sec., Rej {
    if 1
        arg=arg;

    return 1;
}

```

Try an iteration statement [II.2.5] instead of a conditional, and supply an unexpected else clause [II.2.4]:

```

II.2-c5-yes main(arg)
0.26 sec., Acc {
    while(1)
        arg=arg;

    return 1;
}

```

```

II.2-c6-no main(arg)
0.31 sec., Rej {
    while(1)
        arg=arg;
    else
        1;

    return 1;
}

```

To conclude the list of possible errors in the syntax of statements [II.2], consider a valid return statement [II.2.6], a return statement with forgotten semicolon, and a return statement without a value.

```

II.2-d0-yes main(arg)
0.14 sec., Acc {
    return 2+2;
}

```

```

II.2-d1-no main(arg)
0.14 sec., Rej {
    return 2+2
}

```

```

main(arg)
{
    return;
}

```

II.2-d2-no
0.11 sec., Rej

The next two small programs show valid function declarations [II.3.2]:

```

f()
{
    return 1;
}

```

II.3-a0-yes
0.25 sec., Acc

```

main(arg)
{
    return f();
}

```

```

f(x, y)
{
    return 1;
}

```

II.3-a1-yes
0.29 sec., Acc

```

main(arg)
{
    return f(1, 1);
}

```

Using differently ill-formed function headers [II.3.1] leads to incorrect programs:

```

f(x,)
{
    return 1;
}

```

II.3-a2-no
0.21 sec., Rej

```

main(arg)
{
    return 1;
}

```

```

f(x y)
{
    return 1;
}

```

II.3-a3-no
0.21 sec., Rej

```

main(arg)
{
    return 1;
}

```

```

(x)
{
    return 1;
}

```

II.3-a4-no
0.18 sec., Rej

```

main(arg)
{
    return 1;
}

```

II.3-a5-no
0.20 sec., Rej

```
f(,)
{
    return 1;
}

main(arg)
{
    return 1;
}
```

Also, if a valid function header is followed by anything but a compound statement, the program becomes ill-formed:

II.3-a6-no
0.18 sec., Rej

```
f()
    return 1;

main(arg)
{
    return 1;
}
```

II.3-a7-no
0.15 sec., Rej

```
main(arg)
{
    return 1;
}

f()
```

Possible violations of the definition of a program [II.4] as a sequence of functions could be illustrated by the following single statement outside a function:

II-4-a0-no
0.04 sec., Rej

```
return 1;
```

Let us now proceed to the most interesting part of the tests – those covering the additional syntactical constraints [III], most of which cannot be described using context-free grammars.

This well-formed program contains declarations of three functions with different names:

III.1-a0-yes
0.46 sec., Acc

```
function()
{
    return 1;
}

main(arg)
{
    return 1;
}
```

```
function(x, y)
{
    return 1;
}
```

If one symbol is changed in the French name of the third function, making it identical to the name of the first function, the unfortunate fruit of this anglicization is a program that violates the rule of the uniqueness of function declaration [III.1.1]:

```
function()
{
    return 1;
}
```

III.1-a1-no
0.47 sec., Rej

```
main(arg)
{
    return 1;
}
```

```
function(x, y)
{
    return 1;
}
```

Making the headers of these two functions identical is not of much help:

```
function()
{
    return 1;
}
```

III.1-a2-no
0.45 sec., Rej

```
main(arg)
{
    return 1;
}
```

```
function()
{
    return 1;
}
```

Now consider the simplest possible program:

```
main(arg)
{
    return 1;
}
```

III.1-b0-yes
0.14 sec., Acc

Changing the number of arguments of the function `main` makes the program ill-formed [III.1.2]:

```
main(first, second)
{
    return 1;
}
```

III.1-b1-no
0.19 sec., Rej

and so does renaming the main function [III.1.2]:

III.1-b2-no
0.17 sec., Rej

```
principale(arg)
{
    return 1;
}
```

The last way to violate this rule [III.1.2] is to consider a program without any function declarations [II.4], i.e., a possibly empty sequence of whitespace characters:

III.1-c0-no
0.00 sec., Rej

Let us now make some ill-formed function calls [III.2]. First, a well-formed base program:

III.2-a0-yes
0.78 sec., Acc

```
argumentless()
{
    return 1;
}

function(x, y, z)
{
    return x*y*z;
}

main(x)
{
    return argumentless()+function(x-1, x, x+1);
}
```

Changing one letter in the name of the function being called yields a syntactically incorrect program, because now a call to an undeclared function `fonction` is being attempted [III.2.1]:

III.2-a1-no
0.76 sec., Rej

```
argumentless()
{
    return 1;
}

function(x, y, z)
{
    return x*y*z;
}

main(x)
{
    return argumentless()+fonction(x-1, x, x+1);
}
```

Let the function `function` now again be called by its proper name, but move its declaration after its use. The resulting program is ill-formed, because declarations should precede calls [III.2.1]:

```
argumentless()
{
    return 1;
}

main(x)
{
    return argumentless()+function(x-1, x, x+1);
}

function(x, y, z)
{
    return x*y*z;
}
```

III.2-a2-no
0.75 sec., Rej

Calling properly declared functions with an incorrect number of arguments is also not allowed [III.2.2],

```
argumentless()
{
    return 1;
}

function(x, y, z)
{
    return x*y*z;
}
```

III.2-a3-no
0.78 sec., Rej

```
main(x)
{
    return argumentless(x)+function(x-1, x, x+1);
}
```

whether more or less arguments than needed are supplied:

```
argumentless()
{
    return 1;
}

function(x, y, z)
{
    return x*y*z;
}

main(x)
{
    return argumentless()+function(x-1, x+1);
}
```

III.2-a4-no
0.75 sec., Rej

A function can call itself, because in this case its header precedes the reference to it in its body [III.2.1]:

III.2-a5-yes
1.04 sec., Acc

```
argumentless()
{
    return argumentless()+1;
}

function(x, y, z)
{
    return function(x*y*z, 1, argumentless());
}

main(x)
{
    return argumentless()+function(x-1, x, x+1);
}
```

The variable scoping rules [III.3] can be violated in two ways: by defining two variables with the same name and with overlapping scopes [III.4], and by using variables without a prior declaration [III.5].

The first example is a correct program. There are two variables with the same name **fifth**, but their scopes are disjoint:

III.4-a0-yes
1.70 sec., Acc

```
f(first, second, third)
{
    var fourth;

    if(first==second)
    {
        var fifth;
        fifth=second;
        fourth=fifth;
    }
    else
        return third;

    var fifth;
    fifth=fourth;

    return fifth;
}

main(arg)
{
    return f(arg-1, arg, arg+1);
}
```

Let us define another variable named **fourth** inside the **if** statement. The scope of this variable overlaps with the scope of the outer **fourth**, and hence the program becomes ill-formed [III.4]:

```
f(first, second, third)
{
    var fourth;

    if(first==second)
    {
        var fourth, fifth;
        fifth=second;
        fourth=fifth;
    }
    else
        return third;

    var fifth;
    fifth=fourth;

    return fifth;
}
```

III.4-a1-no
1.73 sec., Rej

```
main(arg)
{
    return f(arg-1, arg, arg+1);
}
```

In the next example a variable **second** is defined in a declaration statement [III.3.4], which is itself in the scope of the formal argument **second** [III.3.3]. Thus the program is ill-formed [III.4]:

```
f(first, second, third)
{
    var fourth, second;

    if(first==second)
    {
        var fifth;
        fifth=second;
        fourth=fifth;
    }
    else
        return third;

    var fifth;
    fifth=fourth;

    return fifth;
}

main(arg)
{
    return f(arg-1, arg, arg+1);
}
```

III.4-a2-no
1.73 sec., Rej

The following example declares two more variables called **fifth** [III.3.5], bringing the total count of these variables to 4. However, their scopes are pairwise disjoint, and hence the program remains correct:

III.4-a3-yes
2.12 sec., Acc

```
f(first, second, third)
{
    var fourth;

    if(first==second)
    {
        var fifth;
        fifth=second;
        fourth=fifth;
    }
    else
        return third;

    if(1)
        var fifth;
    else
        while(1)
            var fifth;

    var fifth;
    fifth=fourth;

    return fifth;
}

main(arg)
{
    return f(arg-1, arg, arg+1);
}
```

Let us now leave only one variable with the name `fifth`. The program becomes incorrect, because the identifier `fifth` in the return statement is now outside the scope of the only remaining variable with this name [III.5]:

III.4-a4-no
1.38 sec., Rej

```
f(first, second, third)
{
    var fourth;

    if(first==second)
    {
        var fifth;
        fifth=second;
        fourth=fifth;
    }
    else
        return third;

    return fifth;
}

main(arg)
{
    return f(arg-1, arg, arg+1);
}
```

Following are some simpler examples of multiple variable declarations [III.4]:

```
main(arg)
{
    var longidentifier, x;
    return 1;
}
```

III.4-b0-yes
0.31 sec., Acc

```
main(arg)
{
    var longidentifier, arg;
    return 1;
}
```

III.4-b1-no
0.32 sec., Rej

```
main(arg)
{
    var longidentifier, x, longidentifier;
    return 1;
}
```

III.4-b2-no
0.43 sec., Rej

```
main(arg)
{
    if(1)
        var identifier;
    var identifier;
    return 1;
}
```

III.4-c0-yes
0.40 sec., Acc

```
main(arg)
{
    if(1)
        var identifier, identifier;
    return 1;
}
```

III.4-c1-no
0.36 sec., Rej

```
main(arg)
{
    if(1)
        var identifier, arg;
    var identifier;
    return 1;
}
```

III.4-c2-no
0.43 sec., Rej

```
main(arg)
{
    var identifier;
    if(1)
        var identifier;
    return 1;
}
```

III.4-c3-no
0.45 sec., Rej

```
f(x, y, z)
{
    return 1;
}
```

III.4-d0-yes
0.28 sec., Acc

```
main(arg)
{
    return 1;
}
```

III.4-d1-no
0.25 sec., Rej

```
f(x, y, x)
{
    return 1;
}

main(arg)
{
    return 1;
}
```

The next group of examples refers to the rule of declaration of variables before their use [III.5]. Here is how a simple change of identifiers turns a correct program into ill-formed:

III.5-a0-yes
0.55 sec., Acc

```
main(argument)
{
    var anothervariable;

    return argument+anothervariable;
}
```

III.5-a1-no
0.56 sec., Rej

```
main(argument)
{
    var anothervariable;

    return argument+anothervariable;
}
```

III.5-a2-no
0.56 sec., Rej

```
main(argument)
{
    var anothervariable;

    return argument+uneautrevariable;
}
```

Here is a slightly more complicated case involving scoping rules. The first program is correct:

III.5-b0-yes
0.64 sec., Acc

```
main(argument)
{
    var first;

    if(1)
    {
        var second;
        second=first;
    }

    return first;
}
```

If the expression in the return statement is changed to `second`, the program ceases to be correct, because the scope of this variable is confined to the compound statement inside the conditional statement [III.3.4]:

```
main(argument)
{
    var first;

    if(1)
    {
        var second;
        second=first;
    }

    return second;
}
```

III.5-b1-no
0.64 sec., Rej

Swapping the two statements inside the inner compound statement introduces an error, because the scope of the variable `second` [III.3.4] will begin after its use, and hence the use is ill-formed [III.5]:

```
main(argument)
{
    var first;

    if(1)
    {
        second=first;
        var second;
    }

    return first;
}
```

III.5-b2-no
0.61 sec., Rej

One more special case to test is when the name of a variable coincides with the name of some function. According to the definition of the language, this coincidence is irrelevant, and this is what the grammar demonstrates. In the first well-formed example the name `identifier` is shared by a function and a variable without any side effects:

```
identifier() { return 1; }

main(arg)
{
    var identifier;
    identifier();
    return identifier;
}
```

III.5-c0-yes
0.63 sec., Acc

Modifying a return statement to refer to the function instead of the variable does not spoil the program:

III.5-c1-yes
0.64 sec., Acc

```

identifier() { return 1; }
main(arg)
{
    var identifier;
    identifier();
    return identifier();
}

```

However, an earlier function declaration and a function call together will not atone an undefined variable!

III.5-c2-no
0.43 sec., Rej

```

identifier() { return 1; }
main(arg)
{
    identifier();
    return identifier;
}

```

The last requirement is that every function should return a value [III.6]. The following code manifestly violates this rule:

III.6-a0-no
0.09 sec., Rej

```

main(x)
{
    x=1;
}

```

Here is a well-formed program that can be spoiled in various ways:

III.6-b0-yes
0.56 sec., Acc

```

main(x)
{
    x=x;

    if(2+2<4)
        return 1;
    else if(2+2>4)
    {
        x+x*x;
        return 2;
    }
    else
        return 0;
}

```

If the `return` statement [II.2.6] in any of the three branches of the nested conditional statements [II.2.4] is replaced by anything but a properly returning statement [III.6.1], the program becomes ill-formed [III.6.2]:

```

main(x)
{
    x=x;

    if(2+2<4)
        x=4-(2+2);
    else if(2+2>4)
    {
        x+x*x;
        return 2;
    }
    else
        return 0;
}

```

III.6-b1-no
0.56 sec., Rej

```

main(x)
{
    x=x;

    if(2+2<4)
        while(2+2==4)
            return 1;
    else if(2+2>4)
    {
        x+x*x;
        return 2;
    }
    else
        return 0;
}

```

III.6-b2-no
0.65 sec., Rej

If the `return` statement and the expression-statement [II.2.1] in the properly returning compound statement [III.6.1] in the second branch of the code [III.6-b0-yes] are swapped, then this compound statement ceases to be properly returning, and the program is again incorrect:

```

main(x)
{
    x=x;

    if(2+2<4)
        return 1;
    else if(2+2>4)
    {
        return 2;
        x+x*x;
    }
    else
        return 0;
}

```

III.6-b3-no
0.56 sec., Rej

Taking a simpler well-formed program

```

main(x)
{

```

III.6-c0-yes
0.26 sec., Acc


```

    if(1)
        return 1;
    else
        return 0;
}

```

if the conditional statement loses its else clause [II.2.4], the body of the function `main` is no longer properly returning [III.6]:

III.6-c1-no
0.15 sec., Rej

```

main(x)
{
    if(1)
        return 1;
}

```

The same happens if the conditional statement is replaced with an iteration statement [II.2.5]:

III.6-c2-no
0.17 sec., Rej

```

main(x)
{
    while(1)
        return 1;
}

```

Let us conclude with a reasonable example of a program in this language (reasonable to the extent a program in a language without semantics could be considered reasonable!).

factorial
0.65 sec., Acc

```

factorial(n)
{
    if(n>1)
        return n*factorial(n-1);
    else
        return 1;
}

main(arg)
{
    return factorial(arg);
}

```

If the semantics of the language were defined in the natural way, this program could be said to compute the factorial of a given number.

C Performance

In order to get some impression of the parsing speed as a function of the length of the input, consider a family of programs, each comprised of many simple functions calling each other. The family is parametrized by the number of functions in a program. Here

is the program number ten:

test10
3.14 sec., Acc

```

f00bv(a)
{
    return 1;
}

f01ke(i)
{
    return f00bv(i)*f00bv(i);
}

f02dn(sj)
{
    return f01ke(sj)*f01ke(sj);
}

f03mc(e)
{
    return f02dn(e)*f01ke(e);
}

f04rd(s3)
{
    return f03mc(s3)*f02dn(s3);
}

f05yw(i)
{
    return f04rd(i)*f02dn(i);
}

f06nn(a)
{
    return f05yw(a)*f03mc(a);
}

f07qi(uq)
{
    return f06nn(uq)*f03mc(uq);
}

f08nn(q)
{
    return f07qi(q)*f04rd(q);
}

f09ki(o8)
{
    return f08nn(o8)*f04rd(o8);
}

main(arg)
{
    return f09ki(arg);
}

```

The results are given in Table 1 and plotted in Figure 3; they leave an impression of square time performance.

<i>name</i>	<i>number of symbols</i>	<i>time (sec.)</i>
test01	59	0.33
test03	144	0.88
test05	229	1.48
test10	440	3.14
test15	651	5.05
test20	862	7.23
test25	1073	9.62
test30	1284	12.25
test40	1706	18.17
test50	2128	25.03

<i>name</i>	<i>number of symbols</i>	<i>time (sec.)</i>
test60	2550	32.81
test70	2972	41.59
test80	3394	51.33
test90	3816	62.11
test100	4238	73.98
test115	4901	94.43
test130	5564	117.58
test150	6448	153.22
test175	7553	204.70
test200	8658	265.19

Table 1: Results of performance tests.

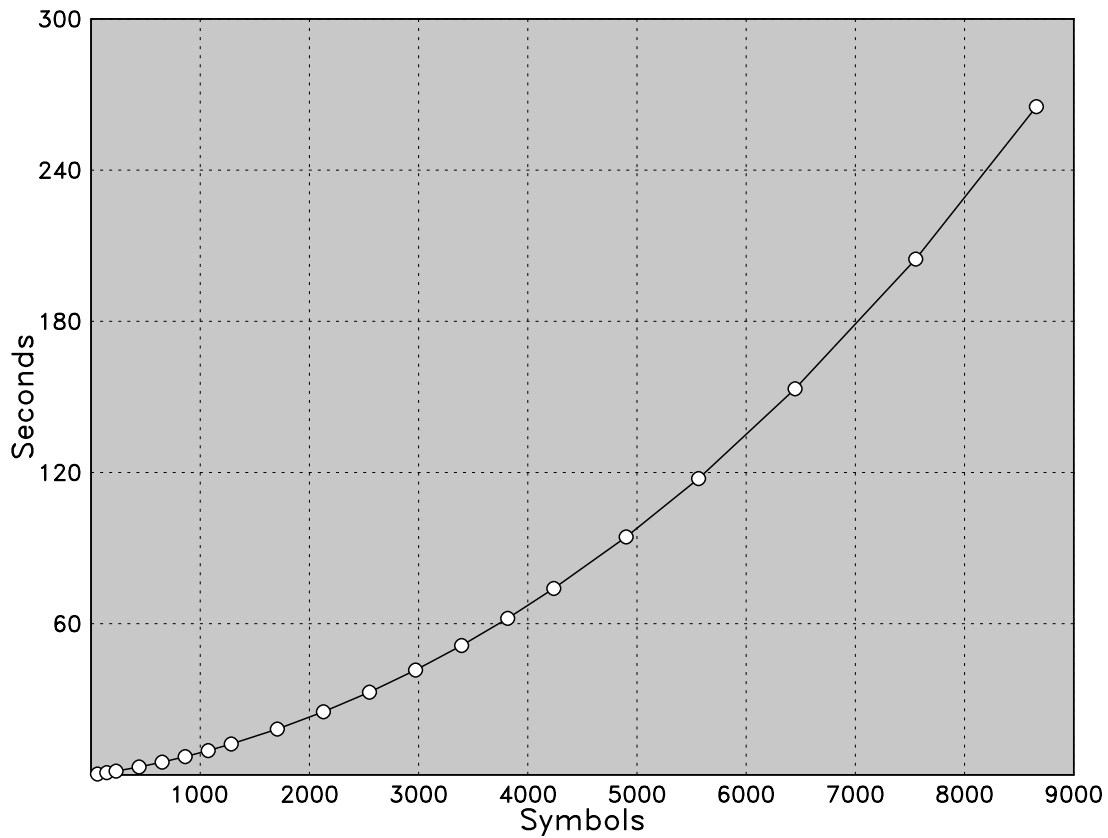


Figure 3: Plotted results of performance tests.