



# Project documentation

## Implementation of the Compiler for the Imperative Language IFJ23

Team xpomsa00  
Variant Hash Table Implementation Using Open Addressing

December 6, 2023

<b>Jakub Pomsár</b>	xpomsa00	25%
Šimon Cagala	xcagal00	25%
Marie Kolářiková	xkolar77	25%
Adrián Ponechal	xponec01	25%

# Contents

<b>1</b>	<b>Teamwork</b>	<b>2</b>
1.1	Work distribution . . . . .	2
1.2	Communication & Git . . . . .	2
1.3	Testing . . . . .	2
<b>2</b>	<b>Lexical analyzer</b>	<b>2</b>
2.1	Implementation . . . . .	2
2.2	Finite-state machine diagram . . . . .	3
<b>3</b>	<b>Symbol Table</b>	<b>4</b>
3.1	Implementation Overview . . . . .	4
3.1.1	Record Structure . . . . .	4
3.2	Hash Functions . . . . .	4
3.3	Key Implementation Aspects . . . . .	4
3.3.1	Item Deletion . . . . .	4
3.3.2	Resizing . . . . .	4
3.3.3	Global and Local Instances . . . . .	4
<b>4</b>	<b>Top-Down Parser</b>	<b>4</b>
4.1	Implementation . . . . .	4
4.1.1	Parser Data . . . . .	4
4.1.2	Token Buffer . . . . .	5
4.2	Analysis . . . . .	5
4.2.1	Syntax . . . . .	5
4.2.2	Semantics . . . . .	5
<b>5</b>	<b>Expression processing</b>	<b>5</b>
5.1	Syntax analysis . . . . .	5
5.1.1	End of expression detection . . . . .	6
5.1.2	Symstack data structure . . . . .	6
5.2	Semantic analysis . . . . .	6
<b>6</b>	<b>Code generator</b>	<b>6</b>
<b>7</b>	<b>Grammar</b>	<b>7</b>
7.1	LL-Grammar . . . . .	7
7.2	LL-Table . . . . .	9

# 1 Teamwork

## 1.1 Work distribution

<b>xpomsa00:</b>	Team leader, symbol table, Top-Down parser, testing, documentation
<b>xcagal00:</b>	Top-Down parser, LL grammar, documentation
<b>xkolar77:</b>	Lexical analyzer, code generator, CI/CD setup, documentation
<b>xponec01:</b>	Precedence analysis, documentation

The workload was evenly distributed, ensuring that each team member contributed equally, and each will receive the same amount of points (25 % ).

## 1.2 Communication & Git

Communication among team members primarily occurred through *Discord* or personal interactions. During the course of the project, we employed the *Git* version control system to manage the evolution of our codebase. *GitLab* served as our remote repository, providing a centralized platform for collaboration and version tracking. The integration of *CI/CD* testing on *GitLab* ensured that changes introduced to the codebase were systematically tested.

## 1.3 Testing

Initially, unit tests were crafted to guarantee the correct behavior of individual modules. As the project progressed, integration tests and end-to-end tests were developed to ensure seamless integration and communication among these modules.

Additionally, trial submissions were employed, with achieving correctness rates of **40%** and **65%**, respectively. Towards the completion of the project, student's tests were incorporated for further validation. As mentioned earlier in 1.2, *CI/CD* testing played a crucial role in verifying the correctness of individual pushed branches.

# 2 Lexical analyzer

The lexical analyzer<sup>1</sup>, implemented as a finite-state machine (2.2), decomposes source code from standard input into a sequence of tokens. Individual tokens are represented by the final states.

## 2.1 Implementation

When the automaton is in the final state and is not able to process the input char, it returns a token to the parser. When the automaton is not in the final state and is not able to process the input char, it returns a lexical error. Each state is represented by a function. The current state is stored in a global variable as a pointer to the corresponding function. Each states sets the function pointer to the next state using the `NEXT_STATE` macro, depending on the loaded char. The token is represented by a structure containing information about the token type, value and flag that indicates the presence of EOL before the processed token. Token value is implemented as a `union` and can contain the name of an identifier, an int, string or double literal value or a boolean `is_nilable`. The token type is implemented using the enum type.

The main loop performs a passage through the automaton by repeated invocation of the function referenced by the variable. Additional global variables, that help remember information that have a lifetime longer than one state, are added in the lexical analyzer. The communication interface with the lexical analyzer is implemented using the function `get_token()`, which has one parameter that represents a pointer to the variable to load the token into. The Lexical analyzer includes counters (`actual_line` and `actual_column`) to tell where the error occurred.

---

<sup>1</sup>Implemented in module `lexical_analyzer.c` with the corresponding header file `lexical_analyzer.h`

## 2.2 Finite-state machine diagram

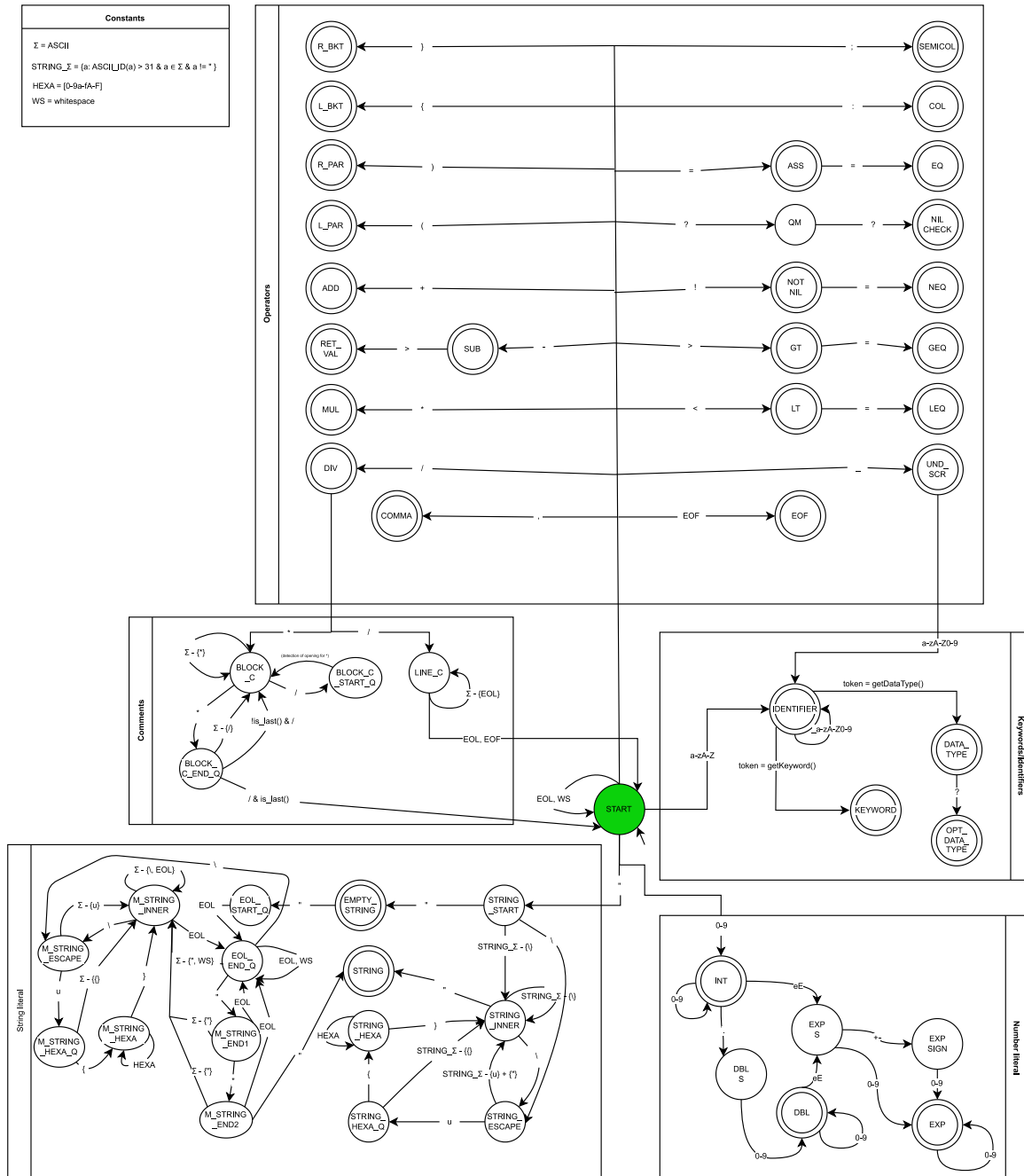


Figure 1: Finite-state machine diagram

## 3 Symbol Table

### 3.1 Implementation Overview

The symbol table<sup>2</sup>, implemented as a hash table using open addressing, dynamically resizes when the load surpasses 65%, handling conflicts through double hashing.

#### 3.1.1 Record Structure

Each record holds essential details: variable/function name, unique identifier, data type, variable mutability status, definition, declaration, nilability, variadic parameters and the return type of functions. The record also points to a linked list preserving parameters.

### 3.2 Hash Functions

Two hash functions are employed: the first (utilized consistently) is *sdbm* variant used in the IJC class [2], and the second, *djb2*, is used for double hashing as introduced in the IAL class [1].

### 3.3 Key Implementation Aspects

#### 3.3.1 Item Deletion

Instead of traditional deletion, the symbol table marks an item's 'activity' attribute as false. Although technically present, the item is considered inactive, creating the appearance of deletion.

#### 3.3.2 Resizing

When the load surpasses **65%**, the symbol table undergoes resizing, creating a new symbol table of size that is two times larger and is a prime number. During this process, every item (excluding inactive items) from the old symbol table is rehashed and transferred to the new one.

#### 3.3.3 Global and Local Instances

A **global symbol table** persists throughout execution, as well as a **stack of local symbol tables**<sup>3</sup> (representing individual scopes), which are *created, stacked, and popped* as needed.

## 4 Top-Down Parser

### 4.1 Implementation

The parser<sup>4</sup> is controlled by the main `parse()` function. It starts off by initializing the `Parser` data structure (4.1.1) and code generator (6) as well as generating the specified built-in functions. Next it utilizes the lexical analyzer to load tokens from `stdin`, which are then stored in the `Parser.buffer` (4.1.2) structure. After all necessary data is initialized, the parser begins two-stage syntax and semantic analysis (4.2), the only exception being expressions, which are parsed using precedence analysis in a separate module<sup>5</sup> (5).

#### 4.1.1 Parser Data

The `Parser` data structure houses all data needed by the parser during analysis, such as: pointers to global symbol table and local symbol table stack (3), the current function, variable or parameter identifiers being processed, current token loaded from the buffer and various counters and flags denoting the state

---

<sup>2</sup>The symbol table is implemented in the files `syntable.c` and `syntable.h`

<sup>3</sup>The stack for local symbol tables is implemented in the files `scope.c` and `scope.h`.

<sup>4</sup>`parser.c` and `parser.h` files

<sup>5</sup>`expression.c` and `expression.h`

the parser is currently in. The data structure is initialized at the start of the `parse()` function and a pointer to it is passed as an argument<sup>6</sup> during recursive descend.

#### 4.1.2 Token Buffer

<sup>7</sup> A doubly linked list of tokens, which are stored for the purpose of parsing the input in two stages (4.2). This structure also allows the parser to backtrack in certain situations, where a simple LL(1) analysis is not sufficient, such as when deciding whether an identifier on the right-hand side of assignment is a variable or called function.

### 4.2 Analysis

For the most part, syntax and semantic analysis happen simultaneously, by recursively calling `Rule` functions which represent the grammar's non-terminals. However, due to the fact that user-defined function calls can be found earlier than their definitions, the input needs to be read twice.

**First Stage** analyzes only function declarations, which entails adding their identifiers, parameters and return type to the global symbol table. Any tokens that are not a part of a function declaration are ignored by the parser using one of the several `skip()` 'Rules'. After EOF is loaded, the `buffer.runner` pointer is reset to the beginning of the list.

**Second Stage** analyzes syntax and semantics of the rest of the input program. At this point the generation of the output code takes place by utilizing the `code_generator` module (6). At this point, when the parser loads tokens related to function declarations, it skips over them to avoid attempts to redefine items already stored in the symbol table.

#### 4.2.1 Syntax

At the centre of syntax analysis lies the `GET_TOKEN` macro, which moves the `buffer.runner` pointer and returns the a token from the list. The types of tokens are then asserted based on rules defined in LL-grammar (7.1). The current token is also used to derive the next `Rule` function to be called. If at any point during the recursive descend the current token does not correspond to the terminals in rule, or the next rule cannot be derived a relevant error code is returned and recursion ends.

#### 4.2.2 Semantics

Semantic analysis is done utilizing mainly symbol table and variables inside `struct Parser` to ensure the input program corresponds with the specifications of IFJ23 language.

Identifiers are checked for attempts of redefinition or usage of undefined ones (which leads to error 3 and 5, depending on the type of identifier).

Function arguments are asserted to be of the same type as the defined parameters and that the correct parameter label is used (or that none is used in the case of a label utilizing the underscore notation). The `Parser.in_func_body` boolean keeps track of whether the parser is currently inside a function, which together with `Parser.return_found` is used for checking the validity of a potential `return` statement. An extra `return` outside of the function body, or a missing one inside leads to error 6.

## 5 Expression processing

### 5.1 Syntax analysis

Whole expression processing is handled by the `expression.c` module with the corresponding header file `expression.h`. The `expr()` function is called whenever an expression is expected by the parser (4). This function accepts a pointer to the `Parser` structure and returns the error code of first error found.

---

<sup>6</sup>e.g. `Rule prog(Parser* p)`

<sup>7</sup>`token_buffer.c` and `token_buffer.h` files

Expressions are processed in 4 states: shift (<), equal shift(=), reduction(>) and error reduction(empty space). The correct state is chosen by using precedence table 5.1.2, where the rows represent the last terminal on stack and the columns represent an incoming terminal in the current expression.

During shift and equal shift operations a symbol is pushed on stack. In the shift state, the handle is set to the closest terminal from the top of the stack. In contrast with the shift state, there is no need to set the handle in equal shift state. The handle is used also in the reduction state.

Expressions are reduced by a suitable rule in reduction state. During this state all symbols are popped from the stack (from top to handle) to an array of symbols. At most 3 symbols can be pushed to this array. These symbols are used to decide what reduction rule should be used for the operation. If there is no rule for these symbols, an error occurs.

### 5.1.1 End of expression detection

End of the expression is detected by the precedence table in error state or in the error handling function in reduction state. If the end is found in error handling functions, recovery is performed (the last pointer is set to the first symbol after the expression and the symbol is removed from stack if necessary).

### 5.1.2 Symstack data structure

A stack implemented in the `symstack.c` module is used for analyzing syntax. It contains inserted terminals and non-terminals of the current expression. This structure is a list of nodes linked from top to bottom. Every node contain data and a pointer to the previous node. The data stores information about a specific expression symbol, like a token, and basic flags that contain information about the corresponding expression symbol. The expression symbol can be either an operand, an operator or a reduced expression.

	*/	+-	??	i	RO	(	)	!	\$
*/	>	>	>	<	>	<	>	<	>
+-	<	>	>	<	>	<	>	<	>
??	<	<	<	<	<	<	>	<	>
i	>	>	>		>		>	>	>
RO	<	<	>	<	=	<	>	<	>
(	<	<	<	<	<	<	=	<	
)	>	>	>		>		>		>
!	>	>	>	>	>	>	>	>	>
\$	<	<	<	<	<	<	<		>

Table 1: Precedence table

Note: *i* - identifier, *RO*- relational operator

## 5.2 Semantic analysis

Expression semantic analysis takes place in the reduction phase of processing the current operation. This is done across several functions with the `process_` prefix. If possible, these functions will trigger implicit conversion of the operands. Implicit conversion is handled by `convert_if_retypeable` function.

The validity of identifiers used in an expression is handled by the `process_operand` function. These functions utilize the `symtable` module to check whether the identifiers were previously defined.

## 6 Code generator

Code generator is implemented in the `code_generator.c` module with the corresponding header file `code_generator.h`. There is no optimizer implemented in the project. Supportive functions are called from the parser or expression modules. The generator utilizes the stack which greatly reduces the number of generated compiler variables. Most constructions are generated directly, the exception being declarations

inside a cycle. Declaration of variables inside loops is overridden on assignments. Inside the cycle, everything except for variable declarations is stored in a buffer, while declarations are written directly to stdout, after the end of the cycle, the entire content of the buffer is written to stdout, this is done to prevent variable redeclarations.

## 7 Grammar

### 7.1 LL-Grammar

1.  $\langle \text{program} \rangle \Rightarrow \langle \text{statement} \rangle \langle \text{program} \rangle$
2.  $\langle \text{program} \rangle \Rightarrow \text{func ID} ( \langle \text{param} - \text{list} \rangle \langle \text{func} - \text{return} - \text{type} \rangle \{ \langle \text{func} - \text{body} \rangle \langle \text{program} \rangle$
3.  $\langle \text{program} \rangle \Rightarrow \text{EOF}$
4.  $\langle \text{statement} \rangle \Rightarrow \langle \text{mutb} \rangle \langle \text{define} \rangle$
5.  $\langle \text{statement} \rangle \Rightarrow \text{ID} \langle \text{id} - \text{type} \rangle$
6.  $\langle \text{statement} \rangle \Rightarrow \text{while EXPRESSION} \{ \langle \text{block} - \text{body} \rangle$
7.  $\langle \text{statement} \rangle \Rightarrow \text{if} \langle \text{cond} - \text{clause} \rangle \{ \langle \text{block} - \text{body} \rangle \text{else} \{ \langle \text{block} - \text{body} \rangle$
8.  $\langle \text{define} \rangle \Rightarrow \text{ID} \langle \text{define} - \text{cont} \rangle$
9.  $\langle \text{define} - \text{cont} \rangle \Rightarrow : \langle \text{type} \rangle \langle \text{opt} - \text{assign} \rangle$
10.  $\langle \text{define} - \text{cont} \rangle \Rightarrow = \langle \text{assign} - \text{exp} \rangle$
11.  $\langle \text{opt} - \text{assign} \rangle \Rightarrow = \langle \text{assign} - \text{exp} \rangle$
12.  $\langle \text{opt} - \text{assign} \rangle \Rightarrow \varepsilon$
13.  $\langle \text{id} - \text{type} \rangle \Rightarrow = \langle \text{assign} - \text{exp} \rangle$
14.  $\langle \text{id} - \text{type} \rangle \Rightarrow ( \langle \text{arg} - \text{list} \rangle$
15.  $\langle \text{cond} - \text{clause} \rangle \Rightarrow \text{let ID}$
16.  $\langle \text{cond} - \text{clause} \rangle \Rightarrow \text{EXPRESSION}$
17.  $\langle \text{arg} - \text{list} \rangle \Rightarrow \langle \text{arg} \rangle \langle \text{arg} - \text{next} \rangle$
18.  $\langle \text{arg} - \text{list} \rangle \Rightarrow )$
19.  $\langle \text{arg} - \text{next} \rangle \Rightarrow , \langle \text{arg} \rangle \langle \text{arg} - \text{next} \rangle$
20.  $\langle \text{arg} - \text{next} \rangle \Rightarrow )$
21.  $\langle \text{arg} \rangle \Rightarrow \text{ID} \langle \text{opt} - \text{arg} \rangle$
22.  $\langle \text{arg} \rangle \Rightarrow \langle \text{literal} \rangle$
23.  $\langle \text{param} - \text{list} \rangle \Rightarrow \langle \text{param} \rangle \langle \text{param} - \text{next} \rangle$
24.  $\langle \text{param} - \text{list} \rangle \Rightarrow )$
25.  $\langle \text{param} - \text{next} \rangle \Rightarrow , \langle \text{param} \rangle \langle \text{param} - \text{next} \rangle$
26.  $\langle \text{param} - \text{next} \rangle \Rightarrow )$
27.  $\langle \text{param} \rangle \Rightarrow \text{LabelID ParamID} : \langle \text{type} \rangle$
28.  $\langle \text{param} \rangle \Rightarrow \_ \text{ParamID} : \langle \text{type} \rangle$



29.  $\langle \text{block} - \text{body} \rangle \Rightarrow \langle \text{statement} \rangle \langle \text{block} - \text{body} \rangle$
30.  $\langle \text{block} - \text{body} \rangle \Rightarrow \}$
31.  $\langle \text{func} - \text{body} \rangle \Rightarrow \langle \text{func} - \text{statement} \rangle \langle \text{func} - \text{body} \rangle$
32.  $\langle \text{func} - \text{body} \rangle \Rightarrow \}$
33.  $\langle \text{func} - \text{statement} \rangle \Rightarrow \langle \text{mutb} \rangle \langle \text{define} \rangle$
34.  $\langle \text{func} - \text{statement} \rangle \Rightarrow \text{ID} \langle \text{id} - \text{type} \rangle$
35.  $\langle \text{func} - \text{statement} \rangle \Rightarrow \text{while EXPRESSION} \{ \langle \text{func} - \text{body} \rangle$
36.  $\langle \text{func} - \text{statement} \rangle \Rightarrow \text{if} \langle \text{cond} - \text{clause} \rangle \{ \langle \text{func} - \text{body} \rangle$
37.  $\langle \text{func} - \text{statement} \rangle \Rightarrow \text{return} \langle \text{opt} - \text{ret} \rangle$
38.  $\langle \text{func} - \text{return} - \text{type} \rangle \Rightarrow \rightarrow \langle \text{type} \rangle$
39.  $\langle \text{func} - \text{return} - \text{type} \rangle \Rightarrow \varepsilon$
40.  $\langle \text{opt} - \text{ret} \rangle \Rightarrow \text{EXPRESSION}$
41.  $\langle \text{opt} - \text{ret} \rangle \Rightarrow \varepsilon$
42.  $\langle \text{opt} - \text{arg} \rangle \Rightarrow : \langle \text{term} \rangle$
43.  $\langle \text{opt} - \text{arg} \rangle \Rightarrow \varepsilon$
44.  $\langle \text{mutb} \rangle \Rightarrow \text{var}$
45.  $\langle \text{mutb} \rangle \Rightarrow \text{let}$
46.  $\langle \text{assign} - \text{exp} \rangle \Rightarrow \text{EXPRESSION}$
47.  $\langle \text{assign} - \text{exp} \rangle \Rightarrow \langle \text{func} - \text{call} \rangle$
48.  $\langle \text{func} - \text{call} \rangle \Rightarrow \text{ID} ( \langle \text{arg} - \text{list} \rangle$
49.  $\langle \text{type} \rangle \Rightarrow \text{Int}$
50.  $\langle \text{type} \rangle \Rightarrow \text{Double}$
51.  $\langle \text{type} \rangle \Rightarrow \text{String}$
52.  $\langle \text{term} \rangle \Rightarrow \text{ID}$
53.  $\langle \text{term} \rangle \Rightarrow \langle \text{literal} \rangle$
54.  $\langle \text{literal} \rangle \Rightarrow \text{INT\_LITERAL}$
55.  $\langle \text{literal} \rangle \Rightarrow \text{DOUBLE\_LITERAL}$
56.  $\langle \text{literal} \rangle \Rightarrow \text{STRING\_LITERAL}$
57.  $\langle \text{literal} \rangle \Rightarrow \text{nil}$

## 7.2 LL-Table

	func	ID	let	var	if	else	while	return	(	)	{	}	EXP	_	:	=	,	->	Int	Double	String	INT_LIT	DBL_LIT	STR_LIT	nil	EOF	ε
<program>	2	1	1	1	1		1																			3	
<statement>		5	4	4	7		6																				
<define>		8																									
<define-cont>															9	10											12
<opt-assign>																11	13										
<id-type>									14																		
<cond-clause>			15										16														
<arg-list>		17							18														17	17	17	17	
<arg-next>									20								19										
<arg>		21																					22	22	22	22	
<param-list>		23							24																		
<param-next>									26								25										
<param>		27																									
<block-body>		29	29	29	29		29					30		28													
<func-body>		31	31	31	31		31	31				32															
<func-statement>		34	33	33	36		35	37																			
<func-return-type>																		38								39	
<opt-ret>													40													41	
<opt-arg>														42												43	
<mult>			45	44																							
<assign-exp>			47																								
<func-call>			48																								
<type>																			49	50	51						
<term>		52																				53	53	53	53		
<literal>																						54	55	56	57		

## References

- [1] Filip Stanis. 008 - djb2 hash. online, September 2019. <https://theartincode.stanis.me/008-djb2>.
- [2] Ozan Yigit. Hash functions. online, September 2003. <http://www.cse.yorku.ca/~oz/hash.html>.