

A COMPILER FOR TINYJAVA - LANGUAGE PROCESSORS

IGNACIO IKER PRADO RUJAS

13th of September 2016

Abstract

In this document we cover a basic overview about the implemented compiler for a language that looks like `Java`, but that in no case achieves all of its potential. We name it `TinyJava`.

A compiler combines different phases in which several analysis are performed, each of them relying on the previous ones. It begins by analyzing the words that form the given code, and it associates to each of them different *lexical units*. Afterwards, it ensures that the sentences are well formed by means of a *context-free grammar*. The next analysis verifies *identifier restrictions* of the variables and performs a *type analysis*. Lastly, it generates *portable code* that can be executed on a *p-machine*.

1. Lexical analyzer

All regarding this analysis is inside the `lexical_analyzer` package. The class named `LexicalAnalyzerTiny` is generated by `JFlex`. The input for `JFlex` is the `grammar.flex` file.

In this first analysis words are explored in order to classify them into different lexical classes, which are symbols with certain associated attributes. These ones are mainly the *lexeme* (meaning the actual word: a *String*), its location inside the code file (line and column) and a certain class represented by an integer.

For instance, it can detect comments that are ignored, but if the word *double* is found, it returns a new lexical unit that will represent this word and its properties. This matching is performed by means of regular expressions, defined in `grammar.flex`.

2. Syntactic analyzer

The code of this part of the compiler is in the `syntactic_analyzer` package. The classes `LexicalClass` and `SyntacticAnalyzerTiny` are generated with `CUP`, thanks to the `LALR.Parser.cup` file.

Working alongside the lexical analyzer discussed above, it checks that the code file is fine in terms of its syntax, which is described by a context-free grammar and summarized as follows:

$$P \rightarrow \text{CLASSLIST } \text{main} \{ \text{STMLIST} \} \\ | \text{main} \{ \text{STMLIST} \}$$

$$\text{CLASSLIST} \rightarrow \text{CLASSLIST } \text{CLASSDEF} | \text{CLASSDEF}$$

$$\text{CLASSDEF} \rightarrow \text{class } \text{CLASSNAME} \{ \text{ATTRIBLIST } \text{METHLIST} \} \\ | \text{class } \text{CLASSNAME} \{ \text{METHLIST } \text{ATTRIBLIST} \}$$

$$\text{ATTRIBLIST} \rightarrow \text{ATTRIBLIST } \text{ATTRIBDECL} ; \\ | \text{ATTRIBDECL} ;$$

$$\text{ATTRIBDECL} \rightarrow \text{public } \text{BASICTYPE } \text{ID} \\ | \text{private } \text{BASICTYPE } \text{ID}$$

$$\text{ATTRIBOBJ} \rightarrow \text{ID} . \text{ID}$$

$$\text{METHLIST} \rightarrow \text{METHLIST } \text{METHOD} \\ | \text{CONSTRUCTOR}$$

$$\text{CONSTRUCTOR} \rightarrow \text{public } \text{CLASSNAME } \text{MPARAMS} \{ \text{STMLIST} \}$$

$$\text{METHOD} \rightarrow \text{public } \text{BASICTYPE } \text{ID } \text{MPARAMS} \{ \text{STMLIST } \text{return } \text{EXPR} ; \} \\ | \text{public } \text{BASICTYPE } \text{ID } \text{MPARAMS} \{ \text{return } \text{EXPR} ; \} \\ | \text{private } \text{BASICTYPE } \text{ID } \text{MPARAMS} \{ \text{STMLIST } \text{return } \text{EXPR} ; \} \\ | \text{private } \text{BASICTYPE } \text{ID } \text{MPARAMS} \{ \text{return } \text{EXPR} ; \} \\ | \text{public void } \text{ID } \text{MPARAMS} \{ \text{STMLIST} \} \\ | \text{private void } \text{ID } \text{MPARAMS} \{ \text{STMLIST} \}$$

$$\text{MPARAMS} \rightarrow (\text{MPARAMSLIST}) \\ | ()$$

$$\text{MPARAMSLIST} \rightarrow \text{MPARAMSLIST} , \text{MPARAM} \\ | \text{MPARAM}$$

$$\text{MPARAM} \rightarrow \text{BASICTYPE } \text{ID}$$

$$\begin{aligned} STMLIST &\rightarrow STMLIST\ STM \\ &\mid STM \end{aligned}$$
$$\begin{aligned} STM &\rightarrow NEWDECL ; \mid DECL ; \mid ASSIGN ; \mid ASSIGNATTRIB ; \\ &\mid ASSIGNARRAY ; \mid WHILE \mid FOR \mid SWITCH \mid IF_ELSE \\ &\mid VOIDMETHOBJ \mid ARRAYDECL ; \end{aligned}$$
$$NEWDECL \rightarrow CLASSNAME\ ID = \text{new}\ CLASSNAME\ (\ EXPRLIST)$$
$$DECL \rightarrow BASICTYPE\ ID \mid BASICTYPE\ ID = EXPR$$
$$ASSIGN \rightarrow ID = EXPR$$
$$ASSIGNATTRIB \rightarrow ATTRIBOBJ = EXPR$$
$$ASSIGNARRAY \rightarrow ARRAYELEM = EXPR$$
$$ARRAYELEM \rightarrow ID\ [\ EXPR \]$$
$$ARRAYDECL \rightarrow BASICTYPE\ [\]\ ID = \text{new}\ BASICTYPE\ [\ INT \]$$
$$BASICTYPE \rightarrow \text{int} \mid \text{bool} \mid \text{double} \mid \text{char} \mid \text{String}$$
$$WHILE \rightarrow \text{while}\ EXPR \{ \ STMLIST \}$$
$$\begin{aligned} IF_ELSE &\rightarrow \text{if}\ EXPR \{ \ STMLIST \} \text{else} \{ \ STMLIST \} \\ &\mid \text{if}\ EXPR \{ \ STMLIST \} \end{aligned}$$
$$FOR \rightarrow \text{for}\ (\ ASSIGN ; EXPR ; ASSIGN \) \{ \ STMLIST \}$$
$$SWITCH \rightarrow \text{switch}\ EXPR \{ \ SWITCHBODY \}$$
$$\begin{aligned} SWITCHBODY &\rightarrow SWITCHBODY\ SWITCHSTM \\ &\mid SWITCHSTM \end{aligned}$$
$$SWITCHSTM \rightarrow \text{case}\ INT \{ \ STMLIST \}$$

$$\begin{aligned}
\text{VOIDMETHOBJ} &\rightarrow \text{ID} . \text{ID} (\text{EXPRLIST}) \mid \text{ID} . \text{ID} () \\
\text{EXPRLIST} &\rightarrow \text{EXPRLIST} , \text{EXPR} \mid \text{EXPR} \\
\text{METHOBJ} &\rightarrow \text{ID} . \text{ID} (\text{EXPRLIST}) \mid \text{ID} . \text{ID} () \\
\text{EXPR} &\rightarrow \text{EXPR} \mid \mid \text{EXPRAND} \mid \text{EXPRAND} \\
\text{EXPRAND} &\rightarrow \text{EXPRAND} \&\& \text{EXPREQ} \mid \text{EXPREQ} \\
\text{EXPREQ} &\rightarrow \text{EXPREQ} == \text{EXPRREL} \mid \text{EXPREQ} != \text{EXPRREL} \mid \text{EXPRREL} \\
\text{EXPRREL} &\rightarrow \text{EXPRREL} < \text{EXPRADD} \mid \text{EXPRREL} > \text{EXPRADD} \\
&\mid \text{EXPRREL} <= \text{EXPRADD} \mid \text{EXPRREL} >= \text{EXPRADD} \mid \text{EXPRADD} \\
\text{EXPRADD} &\rightarrow \text{EXPRADD} + \text{EXPRMULT} \mid \text{EXPRADD} - \text{EXPRMULT} \\
&\mid \text{EXPRMULT} \\
\text{EXPRMULT} &\rightarrow \text{EXPRMULT} * \text{EXPRUNARY} \mid \text{EXPRMULT} / \text{EXPRUNARY} \\
&\mid \text{EXPRMULT} \% \text{EXPRUNARY} \mid \text{EXPRUNARY} \\
\text{EXPRUNARY} &\rightarrow + \text{EXPRUNARY} \mid - \text{EXPRUNARY} \\
&\mid ! \text{EXPRUNARY} \mid \text{EXPRLAST} \\
\text{EXPRLAST} &\rightarrow (\text{EXPR}) \text{ID} \mid \text{INT} \mid \text{DOUBLE} \mid \text{BOOLVAL} \mid \text{CHAR} \mid \text{STRING} \\
&\mid \text{ATTRIBOBJ} \mid \text{METHOBJ} \mid \text{ARRAYELEM}
\end{aligned}$$

Previously on the JFlex file we had defined that, for instance, an identifier of a variable consists of characters starting by a letter, or that a *BOOLVAL* can only be **true** or **false**. Also, as the analysis is performed, it generates an abstract syntax tree. The nodes represent the different analyzed units, such as a node–program whose attributes are the class definitions and the list with the statements of the main program. All of this elements, the statements and the expressions are managed by the subpackage **syntax**.

If the analysis detects any syntactic errors, it will show a message on the console. If not, it will proceed to the next phase.

3. Semantic analyzer

The package named `semantic_analyzer` (alongside `error_handle`) collects this part of the code.

In this part the generated syntax tree is explored in order to detect identifier or type errors, such as a variable that is used but not initialized or trying to add an integer expression with a boolean expression. The class `IdAndTypeChecker` is in charged of this, using `IdManager` which basically contains a stack of tables that gather information about each identifier on a certain block. Every time the code enters on a new block, a new table is pushed onto the stack, and when the block ends it is popped. Depending on the type of identifier, different kind of information is saved on the table, and this information is enclosed on a `IdInfo`.

The main tests regarding types are in the `getType()` method, on the `IdAndTypeChecker` class. See the `Errors` class to get an idea of the type of errors that you can expect it to detect.

4. Code generator

The classes needed for generating the *p-code* are in the `code_generator` package.

On this last phase of the compilation, the given code is translated into code that can be executed on a *p-machine*. The `CodeManager` is in charged of organizing this, and saving the right information. Each statement and expression has its own `code()` method (or `code_L()` or `code_R()`) that tells the `CodeManager` what instructions to add, but it can be tricky because of the jumps on the code.

In order to solve this, we use a queue of queues named `pendingCode` that saves the generated code until the required label is obtained. Other helpful information that is kept is a map that associates to each object the name of its class, a second map that associates to each method the line number where it begins and a third map that associates to each class a map that associates to each attribute an integer value that represents the offset of the attribute viewed as an element of a register:

$$\begin{aligned}\text{objClass} &: \text{Object} \longrightarrow \text{Class}, \\ \text{methLine} &: \text{Method} \longrightarrow \text{Initial line number}, \\ \text{attrOffset} &: \text{Class} \longrightarrow \text{Attribute} \longrightarrow \text{Attribute offset}.\end{aligned}$$

In this way, composing $(\text{objClass} \circ \text{attrOffset})$ we obtain for each object the offset of its attributes, which is needed for the code generation.

The generated code is forwarded to a file named `output.txt` inside the `examples` folder, from where the input is taken.