
Project part 2: Parser

Subject :
introduction to
language theory and
compiling

TRANGANIDAS Orestis
BOUGLAM Sara

Contents

1	Introduction.....	3
2	Background.....	4
2.1	Parser.....	4
2.2	Grammar.....	4
2.2.1	Definition.....	4
2.2.2	Chomsky hierarchy.....	4
2.3	Grammar derivation:.....	5
2.4	First ^k and follow ^k	5
2.4.1	First ^k	5
2.4.2	follow ^k	5
2.5	LL(1) parser.....	5
2.5.1	Unproductive and inaccessible rules removal.....	5
2.5.2	Ambiguity removal.....	5
2.5.3	Left recursion removal.....	6
2.5.4	Left prefixes removal.....	6
2.5.5	Action table.....	6
2.5.6	Recursive descent.....	6
3	Implementation.....	7
3.1	Overview.....	7
3.2	Phase 1.....	7
3.2.1	Unproductive characters.....	7
3.2.2	Unreachable characters.....	8
3.2.3	Ambiguous grammar.....	8
3.2.4	Left Factory.....	9
3.2.5	Left recursion.....	10
3.2.6	LL(1) grammar.....	10
3.2.7	First ¹ and Follow ¹	12
3.2.8	Action Table.....	12
3.3	Phase 2.....	14
3.3.1	Error message.....	14

3.3.2	Extra feature.....	15
3.3.3	Main class	15
4	Conclusion	16
5	Bibliography	16

1 Introduction

The compiling process goes through several phases, starting by the lexical analysis where the deterministic finite automaton is introduced as a tool to model and execute the process, resulting in the generation of a set of finite lexical units representing the whole program; afterwards, the syntactic analysis comes in second place, which treats the set of finite symbols resulted in the prior phase.

Deterministic finite automaton might be a suitable tool to tackle the parser generation, however, it is flawed and has numerous limitations, this can be noticed in the L_0 ; such a simple language that is not regular, therefore, the concept of grammar is brought up, which fill the gaps of the deterministic finite automaton .

Formally, grammar simply describes how the string is formed from a language's alphabets that match the language's syntax.

In this report, parser generation process will be addressed. Initially, $LL(1)$ grammar accompanied with the symbol table are generated ,then , as a final step the parser is implemented programmatically using recursive descent, in addition to create the parse tree regarding the inputted Fortran program which validates its syntactic correctness.

2 Background

2.1 Parser

The parser is responsible of performing the syntactic analysis of a program, where it receives as an input the set of tokens generated by the combination of the scanner and the screener and output the syntax or the parse tree that represent the structure of the whole program which would be accessed by the compiler in a later phase, two type of parser can be distinguished; *top down parser* and *bottom up parser*.

When the syntax error occurs in a program, the parser is supposed to localized and report it, optionally, diagnose it and correct it (1).

2.2 Grammar

2.2.1 Definition

The syntax analysis is characterized by the grammar, which consists of the elementary component of a program. (1)

Grammar is a quadruplet (2) $G = (V, T, P, S)$

V is a finite set of variables

T is a finite set of terminals

P is a finite set of production rules $\alpha \rightarrow B$ such that :

- $\alpha \in (V \cup T)^* V (V \cup T)^*$
- $B \in (V \cup T)^*$

2.2.2 Chomsky hierarchy

Grammars are divided into 4 distinguishable classes

- *Class 0*: includes all the classes(1,2,3)
- *Class 1*: Context Sensitive Grammar
$$\alpha \rightarrow B \in P \mid \alpha = S \wedge B = \varepsilon \vee |\alpha| \leq |B|/\{S\}$$
- *Class 2*: Regular Grammar
$$\alpha \rightarrow B \in P \mid \alpha \in V \wedge |\alpha| = 1$$
- *Class3*: Context free grammar: is comprised of :
 1. *Left regular grammar* $\alpha \rightarrow B \mid \alpha \in P \wedge (B \in T^* \vee B \in T^*.V)$
 2. *Right regular grammar* $\alpha \rightarrow B \mid \alpha \in P \wedge (B \in T^* \vee B \in V.T^*)$

The syntactic phase of a compiler is written using Context free grammar.

A context free language (CFL) is defined such that $L(G) = L$ where G is context free grammar (2).

2.3 Grammar derivation:

The derivation of a string for a specific grammar is the process of generating a sequence of grammar rules application that does the start symbol transformation into the string (3). There exist two types of derivations, left most and right most one (2):

For a derivation $w, S w' \Rightarrow w \alpha w'$:

- Left most derivation requires that $w \in T^*$
- Right most derivation requires that $w' \in T^*$

In order to prove that given word belongs to the language of a grammar, *derivation tree* can be used, note that the tree is completed when leaves only contain terminals (2)

2.4 First^k and follow^k

2.4.1 First^k

$$\text{First}^k(a) = \{ w \in T^* \mid a \Rightarrow^* wx \wedge (|w| = k \vee |w| < k \text{ and } x = \epsilon) \}$$

2.4.2 follow^k

$$\text{Follow}^k(a) = \{ w \in T^* \mid \{ \exists \beta, \gamma \mid \beta \Rightarrow^* \beta a \gamma \wedge w \in \text{First}^k(\gamma) \} \}$$

2.5 LL(1) parser

LL refers to Left scanning Left parsing which means the input string is scanned from the left to the right. $LL(K)$ grammar uses K *look-ahead* which defines the generation of k *firsts* and *follows*, therefore, $LL(1)$ grammars can be defined as an $LL(k)$ grammar where $K = 1$

However, the grammar should adhere to a set of rules in order not to confuse the parser later on, for this, certain grammar transformation are required initially (2):

2.5.1 Useless rules removal

For a given grammar $G = (V, T, P, S)$, useless rules removal consists of eliminating the following the follow component:

- Unproductive variable: $\{ \text{there is no } w \in T^* \mid A \Rightarrow_G^* w \}$
- Unreachable symbols: $\{ \text{for } X \in V \cup T, \nexists S \Rightarrow_G^* \alpha 1 X \alpha 2 \}$

2.5.2 Ambiguity removal

Grammar is considered ambiguous In the cane when several derivative tree are extracted for the same word.

Ambiguity can be removed through adjusting the production rules such that, they respect a certain order (2).

2.5.3 Left recursion removal

Left recursion can be illustrated such that:

$$S \rightarrow Sa$$

$$S \rightarrow c$$

One can see that: $\text{First}(c) \in \text{first}(S) \subseteq \text{first}(Sa)$, which violate the $LL(1)$ grammar, this can be fixed by turning the left recursion into right recursion (2).

2.5.4 Left prefixes removal

The following grammar is left prefixed:

$$A \rightarrow aB$$

$$A \rightarrow aC$$

This grammar can be adjusted through factoring the common prefixes.

2.5.5 Action table

After transforming the grammar into $LL(1)$ one, the action table is considered as the core of the parser; it is obtained through the result of generating the follow1 and first1 from the transformed grammar, it is a two dimensional table M (2) :

- The lines are indexed by the elements from $T \cup V$
- Rows of M are indexed by a set of T
- Cells contain the action that the parser must perform

2.5.6 Recursive descent

Recursive descent is a type of top down parser which relies on the action table obtained as well as the $LL(1)$ grammar (2):

For the input, it receives the following:

- $LL(1) CFG = (P, T, V, S)$ accompanied with the action table and an input word
- $w = w_1, w_2, \dots \in T^*$

And outputs :*True* iff $w \in L(G)$, the sequence of rules number is printed in a left most derivation

3 Implementation

3.1 Overview

The process of generating the parser was divided into two main phases, the first phase required manual intervention which consisted of transforming the **Fortran** grammar to **LL(1)** and action table generation, then the second phase relied on the results obtained from the first phase, eventually, implementing the recursive descent using java programming language, then generating the parse tree.

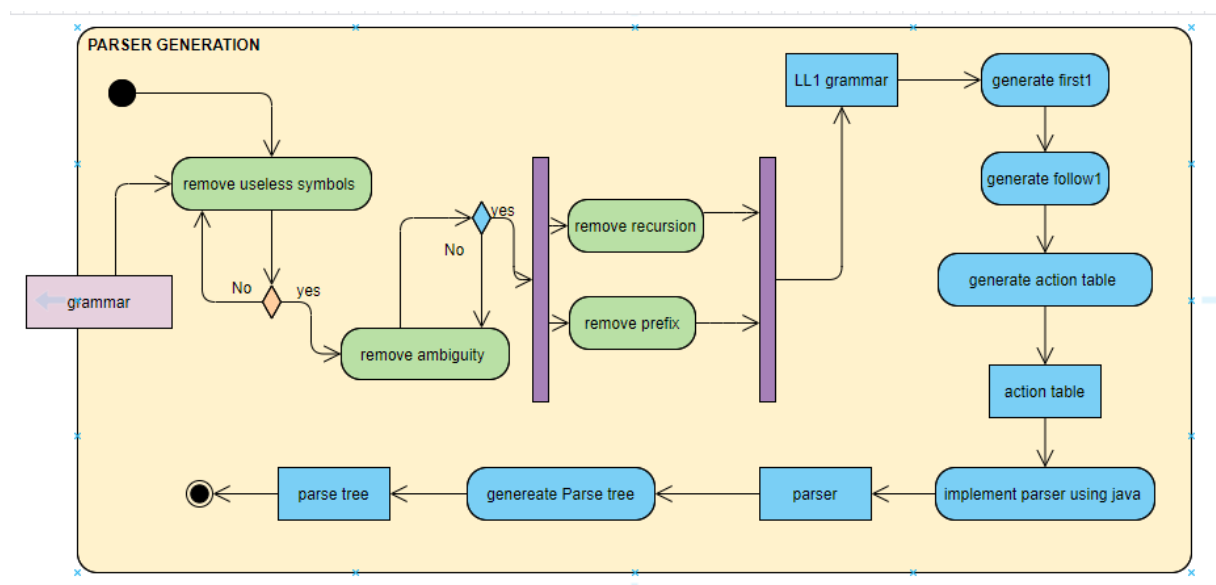


Figure 1: An activity diagram¹ that summarize the overall process using visual paradigm web app²

3.2 Phase 1

The subsequent sections explain the whole process followed in order to generate the action table

3.2.1 Unproductive characters

Considering that the set of terminals corresponds to the units contained in the Symbol Class

I	V_i
0	ϕ
1	{<Read>}
2	{<Read>, <Print> }
3	{<Read>, <Print>, <While>}
4	{ <Read>, <Print>, <While>, <Comp> }

¹ **Activity Diagram** : is a flow chart that represents the flow from one activity to another in a system, this activity can be described as an operation (6)

² **Online visual paradigm**: is a web app used to generate different template and diagrams
<https://online.visual-paradigm.com/>

5	{<Read>, <Print>, <While>, <Comp> <Cond>}
6	{<Read>, <Print>, <While>, <Comp> <Cond>, <If>}
7	{<Read>, <Print>, <While>, <Comp> <Cond>, <If>, <Op>}
8	{<Read>, <Print>, <While>, <Comp> <Cond>, <If>, <Op>, <ExprArith>}
9	{<Read>, <Print>, <While>, <Comp> <Cond> , <If> , <Op> , <ExprArith> , <Assign>}
10	{<Read>, <Print>, <While>, <Comp> <Cond>, <If>, <Op>, <ExprArith> , <Assign> , <Instruction>}
12	{<Read>, <Print>, <While>, <Comp> <Cond>, <If> , <Op>, <ExprArith>, <Assign> , <Instruction>, <Code>}
13	{<Read>, <Print>, <While>, <Comp> <Cond>, <If>, <Op>, <ExprArith>, <Assign> , <Instruction> , <Code>, <Program>}

Observation: There are no unproductive rules in the given grammar

3.2.2 Unreachable characters

I	V _i
0	{<Program>}
1	{<Program> , <Code>}
2	{<Program>, <Instruction>}
3	{<Program>, <Instruction>, <Assign> , <While> , <Print>, <Read>}
4	{<Program>, <Instruction>, <Assign>, <While>, <Print>, <Read>, <ExprArith>}
5	{<Program>, <Instruction>, <Assign>, <While>, <Print>, <Read>, <ExprArith>, <Op>}
6	{<Program>, <Instruction>, <Assign>, <While>, <Print>, <Read>, <ExprArith>, <Op>, <If>}
7	{<Program>, <Instruction>, <Assign>, <While>, <Print>, <Read>, <ExprArith>, <Op>, <If>, <Cond>}
8	{<Program>, <Instruction>, <Assign>, <While>, <Print>, <Read>, <ExprArith>, <Op>, <If>, <Cond>, <Code>}
9	{<Program>, <Instruction>, <Assign>, <While>, <Print>, <Read>, <ExprArith>, <Op>, <If>, <Cond>, <Code>, <Program>}

Observation: There are no unreachable rules in the given grammar

To sum up, there are no useless rules.

3.2.3 Ambiguous grammar

In the given grammar, there is ambiguity regarding the addition and multiplication operation which need to be adjusted according to their order

```

<ExprArith> → [VarName]
            → [Number]
            → ( <ExprArith> )
            → - <ExprArith>

```

```

→ <ExprArith> <Op> <ExprArith>

<Op> → +
      → -
      → *
      → /

```

After the adjustment by prioritizing the multiplication and division, the above grammar is obtained

```

<ExprArith> -> <ExprArith> + <Multiplication>
             -> <ExprArith> - <Multiplication>
             -> <Multiplication>
<Multiplication> -> <Multiplication> * <Bracket>
                  -> <Multiplication> / <Bracket>
                  -> <Bracket>
<Bracket> -> (<ExprArith>)
           -> <Var>

<Var> -> [VarName]
        → [Number]
        → - <Var>

```

3.2.4 Left Factory

Now that ambiguity and useless rules have been eliminated, there is a noticeable prefix redundancy in the following rules:

```

<If> → IF (<Cond>) THEN [EndLine] <Code> ENDIF
<If> → IF (<Cond>) THEN [EndLine] <Code> ELSE [EndLine] <Code> ENDIF

```

After the removal of the left factory we get

```

<If> → IF (<Cond>) THEN [EndLine] <Code> <If''>
<If''> -> ENDIF
<If''> → ELSE [EndLine] <Code> ENDIF

```

3.2.5 Left recursion

Left recursion is noticed in the rule bellow:

```
<ExprArith> -> < ExprArith > + <Multiplication>
              -> < ExprArith > - <Multiplication>
              -> <Multiplication>
<Multiplication> -> <Multiplication> * <Braquet>
                  -> <Multiplication> / <Braquet >
                  -> <Braquet>
```

After the removal of left recursion:

```
< ExprArith > -> <ExprArith'> < ExprArith''>
< ExprArith''> -> + <Multiplication>< ExprArith''>
               -> - <Multiplication>< ExprArith''>
               ->  $\epsilon$ 
<ExprArith'> -> <Multiplication>

<Multiplication> -> <Multiplication'> <Multiplication''>
<Multiplication''> -> * <Braquet> <Multiplication''>
                  -> / <Braquet> <Multiplication''>
                  ->  $\epsilon$ 
<Multiplication'>-> <Bracket>
```

3.2.6 LL(1) grammar

In the end, the following *LL*(1) grammar is obtained

```
[1] <S>   -> <Program>$
[2] <Program> → BEGINPROG [ProgName] [EndLine] <Code> ENDPROG
[3] <Code> → <Instruction> [EndLine] <Code>
[4]       →  $\epsilon$ 
[5] <Instruction> → <Assign>
```

```

[6]          → <If>
[7]          → <While>
[8]          → <Print>
[9]          → <Read>
[10] <Assign> → [VarName] := <ExprArith>
[11] <ExprArith > -> <ExprArith'> < ExprArith''>
[12] <ExprArith''> -> + <Multiplication>< ExprArith''>
[13]          -> - <Multiplication>< ExprArith''>
[14]          ->  $\varepsilon$ 
[15] <ExprArith'> -> <Multiplication>
[16] <Multiplication> -> <Multiplication'> <Multiplication''>
[17] <Multiplication''> -> * <Braquet> <Multiplication''>
[18]          -> / <Braquet> <Multiplication''>
[19]          ->  $\varepsilon$ 
[20] <Multiplication'> -> <Bracket>
[21] <Bracket> -> (ExprArith)
[22]          -> <Var>
[23] <Var> -> [VarName]
[24]          → [Number]
[25]          → - <Var>
[26] <If> → IF (<Cond>) THEN [EndLine] <Code> <If''>
[27] <If''> -> ENDIF
[28] <If''> → ELSE [EndLine] <Code> ENDIF
[29] <Cond> → <ExprArith> <Comp> <ExprArith>
[30] <Comp> → =
[31]          → >
[32] <While> → WHILE (<Cond>) DO [EndLine] <Code> ENDWHILE
[33] <Print> → PRINT([VarName])
[34] <Read> → READ([VarName])

```

3.2.7 First¹ and Follow¹

The last step before the action table generation is the *First*¹ and *Follow*¹ extraction:

Symbol	First ¹	Follow ¹
<S>	BEGINPROG	
<Program>	BEGINPROG	\$
<Code>	VarName IF WHILE PRINT READ ϵ	ENDPROG ENDIF ELSE ENDWHILE
<Instruction>	VarName IF WHILE PRINT READ	ENDLINE
<Assign>	VarName	ENDLINE
<ExprArith>	(VarName Number -	ENDLINE)
<ExprArith''>	+ - ϵ	ENDLINE)
<ExprArith'>	(VarName Number -	+ -
<Multiplication>	(VarName Number -	+ -
<Multiplication''>	/ * ϵ	+ -
<Multiplication'>	(VarName Number -	/ *
<Bracket>	(VarName Number -	/ *
<Var>	VarName Number -	/ *
<If>	IF	ENDLINE
<If''>	ENDIF, ELSE	ENDLINE
<Cond>	(VarName Number -)
<Comp>	= >	(VarName Number -
<While>	WHILE	ENDLINE
<Print>	PRINT	ENDLINE
<Read>	READ	ENDLINE

3.2.8 Action Table

	BE GI NP RO G	PR OG NA ME	()	+	-	*	/	=	>	I F	E L S E	E N D I F	W H I L E	EN DW HI LE	VA RN AM E	N u m b e r	P R I N T	R E A D	EN DP RO G	EN DL IN E
<S>	1																				
<Prog ram>	2																				
<CODE >											3	4	4	3	4	3		3	3	4	

<Instruction>											6			7			5		8	9		
<Assign>																	10					
<Expr Arith>			1 1			1 1											11	1 1				
<Expr Arith''>				1 4	1 2	1 3																14
<Expr Arith'>			1 5			1 5											15	1 5				
<Multiplication>			1 6			1 6											16	1 6				
<Multiplication''>					1 9	1 9	1 7	1 8														
<Multiplication'>			2 0			2 0											20	2 0				
<Bracket>			2 1			2 2											22	2 2				
<Var>						2 5											23	2 4				
<If>										2 6												
<If''>										2 7	2 8											
<Cond>			2 9			2 9											29	2 9				
<Comp>								3 0	3 1													
<While>														3 2								
<Print>																			3 3			
<Read>																				3 4		

3.3 Phase 2

This phase consisted of implementing the recursive descent using Java programming language; the whole idea behind this mechanism consists of generating method for each production rule, and each method is contained into another which receives as a parameter the parent node of the tree and optionally the previously generated token, both of these parameters are used for generating the parse tree.

```
71     private static void program(ParseTree parent) throws java.io.IOException{
72         //Input the current rule in the list of rules that were used
73         if(verbose){
74             rulesText += "[2] <Program> -> BEGINPROG [ProgName] [EndLine] <Code> ENDPROG\n";
75         }else{
76             rulesText += "2 ";
77         }
78         //Create a new Node that represents the rule
79         ParseTree tree = new ParseTree(new Symbol(Labels.PROGRAM));
80         parent.addChild(tree);
81         Symbol token = analyzer.nextToken();
82         //Check if the next token matches the one expected by the rule
83         if(token.getType() == LexicalUnit.BEGINPROG){
84             tree.addChild(new ParseTree(token));
```

Figure 2 Excerpt from the program implemented in the *parser* class

Observation:

In the excerpt above, the *< program >* production rule is displayed, receives as a parameter the parent node which is also the start of the tree.

When the verbose variable (line 73) is set to true, extra information about the parsing would be provided, which consists of displaying the whole production rule.

3.3.1 Error message

Whenever a syntactic error occurs in the input file, an error message is displayed to the standard output stream, which precise the error, the expected token and the exact position (line and column).

```
849     private static void errorMessage(String expected, String received, int line, int column){
850         error = true;
851         errorText = "Error at line "+line+" at column "+column+" : Expected " +expected+ " and instead received "+received;
852         System.out.println(errorText);
853     }
```

Figure 3 : The error message implemented

3.3.2 Extra feature

Additionally, another feature has been added to the *Parser* class that is the variable recording which consists of storing the program variables in a data structure, for later use.

This feature has been implemented using the *map*³ data structure; more precisely the *SortedMap*⁴, the *var* argument is simply map key which refers the variable name and the place *arguments* represents the line where the variable appeared

```
863     private static void record(String var, int place){
864         if(! variables.containsKey(var)){
865             variables.put(var, place);
866         }
867     }
```

Figure 4: Variable recording method

3.3.3 Main class

The main class simply, put all the pieces together: receiving the following arguments from the command line:

1. *-v* stands for the verbose (whether to display extra information about the production rules)
2. *-wt* stands for the file to display the parse tree, *.tex* extension is expected, else, an error is returned
3. *Fortran* program file (*.fs* extension is expected)

Afterward, parsing the program and outputting the parse tree into the latex file.

³ **Map**, java util package's class data structure that consists of mapping a specific value to a key, where each key can map at most one value (4)

⁴ **Sorted map** implements the *map* interface which provide a totally ordering of its keys according to a natural order or by comparator (5)

4 Conclusion

To sum up, parser generation relied on two main steps, the first step was reformulating the grammar into an appropriate one that can be understood by the parser which is LL(1) grammar, then, filling the action table through the obtained first 1 and follow1, the second step consisted of implementing the recursive descent using java programming language, which included some additional features, namely : recognizing syntactic errors and their exact position, generating a parse tree for the programs which language is syntactically correct and recording variable appearance.

Bibliography

1. **Reinhard Wilhelm, Helmut Seidl, Sebastian Hack.** *Compiler Design* .
2. **GILLES GEERAERTS, GUILLERMO A. PÉREZ.** *Introduction to language theory and compilig*.
3. Context-free grammar. *Wikipedia* . [Online] https://en.wikipedia.org/wiki/Context-free_grammar#Derivations_and_syntax_trees.
4. Interface Map. *Java docs* . [Online] <https://docs.oracle.com/javase/7/docs/api/java/util/Map.html>.
5. SortedMap. *Java Doc*. [Online] <https://docs.oracle.com/javase/7/docs/api/java/util/SortedMap.html>.
6. tutorials point. *UML - Activity Diagrams*. [Online] https://www.tutorialspoint.com/uml/uml_activity_diagram.htm.
7. *Introduction to language theory and compillig* .