

Faculty of Engineering, Alexandria University, Computer and Systems Engineering Department

Phase 2 Compiler "Parser Generator"

Names:

Yomna Ahmed 71

Nada Ahmed 66

Marwan Tarek 61

Mohamed Elsabbagh 52

A description of the used data structures:

- 1. We depend mainly on HashMap data structure.
 - In the class ParseGenerator, we use two hashmaps:
 HashMap<String,NonTerminal> non_terminal_map
 HashMap<String,Terminal> terminal_map
 to hold the non terminals and the terminals of the
 grammar with their key is the name of the terminal or
 non terminal and the value is a reference to it.
- 2. We also use ArrayList<> data structure.
 - In the NonTerminal class, each non terminal holds its production rules: ArrayList<ArrayList<Component>> where the component is either terminal or non terminal and each production rule is ArrayList<Component>.
 - In the **Pair** class to represent each production where the non terminal is mentioned in the grammar.
 - In the ParseGenerator class, we use two dummy arraylists

ArrayList<Component> dummy_for_epsilon
ArrayList<Component> dummy_for_synch
To represent the synch and the epsilon production in the parse table.

And ArrayList<Component>[][] for the parse table.

- 3. We also use Stack data structure.
 - In the process_input() function, Stack is used to put
 the productions in the stack from the reverse direction
 ,for example (S → Aab if we need to replace S by this
 production in the stack we put it as the following b a A).

Explanation of algorithms and techniques used:

Parse the grammar file:

We read the grammar file line by line to get the production rules and save the terminals and non terminals of the grammar each into its corresponding map. With each non terminal having its productions saved as an arraylist of terminals and non terminals. And also it has the productions, in which it's referred to in toto, to help us compute the follow set of this non terminal.

• Eliminating left recursion:

We iterate over the non terminal map and eliminate the left recursion of each non terminal after we substitute with each preceding non terminal mentioned in the current non terminal's production rules.

To Eliminate the left recursion we use this method:

$$A \rightarrow A \alpha_1 \mid ... \mid A \alpha_m \mid \beta_1 \mid ... \mid \beta_n$$
 where $\beta_1 ... \beta_n$ do not start with A
 \downarrow eliminate immediate left recursion

$$A \to \beta_1 \, A^{\cdot} \, | \, ... \, | \, \beta_n \, A^{\cdot}$$

$$A \rightarrow \alpha_1 A' \mid ... \mid \alpha_m A' \mid \epsilon$$
 an equivalent grammar

• Left factoring:

We iterate over the non terminal map and eliminate the left factoring of each non terminal by check if a non terminal has **two** or more production rules with a common non-empty prefix

$$A \rightarrow \; \alpha\beta_1 \, | \, ... \, | \, \alpha\beta_n \; | \; \; \gamma_1 \, | \, ... \; | \, \gamma_m$$

And convert it to

$$A \to \alpha A' | \gamma_1 | \dots | \gamma_m$$

$$A' \to \beta_1 | \dots | \beta_n$$

• Compute first set:

We iterate over the non terminal map and for each nonterminal check its productions:

- If X a where 'a' is terminal.
 - 'a' is in FIRST(X).
- If X epsilon
 - epsilon is in FIRST(X).
- If X Y where 'Y' is nonterminal.
 - FIRST(Y) is in FIRST(X).
- If X Y1 Y2 Y3 Yn.
 - If epsilon is in FIRST(Yj) where j = 1,2,...,i-1 then FRIST(Yi) is in FIRST(X).
 - If epsilon is in FIRST(Yj) where j = 1,2,...,n then epsilon is in FIRST(X).

Compute follow set:

First, If S is the start symbol --> \$ is in FOLLOW(S). Then, we keep track of productions that a non terminal was mentioned in for each non terminal in the array list of pairs where pair has a source non terminal (the non terminal on LHS) and array list of components that represents the production and apply rules:

- 1. if A --> aBC is a production rule --> everything in FIRST(C) is FOLLOW(B) except epsilon.
- If (A --> aB is a production rule) or (A --> aBC is a production rule and epsilon is in FIRST(C))
 everything in FOLLOW(A) is in FOLLOW(B).

Construct parse table:

Since we have a map containing all non terminals and another one for terminals, we create the parsing table as a 2D array of productions. We map each non terminal to a row and map each terminal to a column in this array. Then we start filling entries in parsing table by iterating over the nonterminal map in this way:

- 1. we get the corresponding row for nonterminal x.
- 2. for each production:
 - 1) If the first component is a terminal, we get the column corresponding to it and put the production in this cell.
 - 2) If the first component is a nonterminal, we get columns corresponding to the first set of this nonterminal and put the production in these cells.
 - 3) If nonterminal x has epsilon production, we get columns corresponding to terminals in its follow set and put dummy for epsilon in these cells.

4) if nonterminal x don't have epsilon production, we get columns corresponding to terminals in its follow set and put dummy for sync in these cells (if these cells were empty).

Parse the input program:

This part is implemented by the process_input function at first stack is created and push in it "\$" sign and the starting Non terminal symbol then loop until stack is empty and inside this loop we check the following:

- If the top of the stack is Terminal then we have to paths:
 - 1) topOfstack name the same as that of the next token then pop from the stack
 - 2) Else Report Error and in the next loop don't take the next token
- Else when the top of the stack is non terminal then we get the production that fit with the token and the Terminal/ Non terminal in the stack and here there are 3 cases, if after looking in the parsing table we find that:
 - 1) there is "sync" then:
 - Report an Error
 - pop the stack and don't take the next token
 - 2) if it is Empty then:
 - Report an Error
 - get next token
 - 3) Otherwise, pop the stack and take the productions of the component pushing from last to first .

Input files:

```
input_RegEx.txt - Notepad
                                                                   input_grammar.txt - Notepad
File Edit Format View Help
                                                                   File Edit Format View Help
letter = a-z | A-Z
                                                                   # METHOD_BODY = STATEMENT_LIST
digit = 0 - 9
                                                                   # STATEMENT LIST = STATEMENT | STATEMENT LIST STATEMENT
id: letter (letter|digit)*
                                                                   # STATEMENT = DECLARATION
digits = digit+
                                                                   | IF
                                                                   WHILE
{boolean int float}
num: digit+ | digit+ . digits ( \L | E digits)
relop: \=\= | !\= | > | >\= | < | <\=
                                                                   ASSIGNMENT
                                                                   # DECLARATION = PRIMITIVE_TYPE ' id' ';'
                                                                  # PRIMITIVE_TYPE = 'int' | 'float'
# IF = 'if' '(' EXPRESSION ')' '{' STATEMENT '}' 'else' '{' STATEMENT '}'
# WHILE = 'while' '(' EXPRESSION ')' '{' STATEMENT '}'
# ASSIGNMENT = 'id' '=' EXPRESSION ';'
assign: \=
{ if else while }
[; , \( \) { }]
addop: \+ | -
mulop: \* | /
                                                                   # EXPRESSION = SIMPLE_EXPRESSION
                                                                   | SIMPLE EXPRESSION 'relop' SIMPLE EXPRESSION
                                                                   # SIMPLE_EXPRESSION = TERM | SIGN TERM | SIMPLE_EXPRESSION 'addop' TERM
                                                                   # TERM = FACTOR | TERM 'mulop' FACTOR
# FACTOR = 'id' | 'num' | '(' EXPRESSION ')'
                                                                   # SIGN = '+' | '-'
```

Test program:

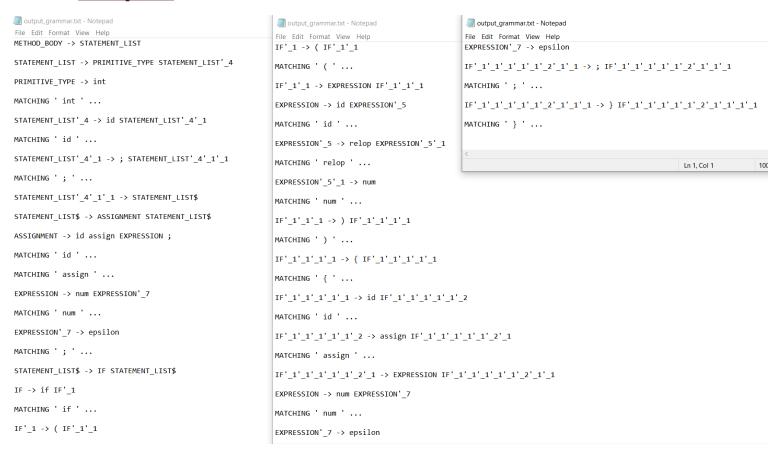
```
int x;
test.txt - Notepad

File Edit Format View Help

int x;
```

```
µnc x;
x = 5;
if (x > 2)
{
x = 0;
}
```

Output:



The resultant stream of tokens for the example test program:

```
output_tokens - Notepad

File Edit Format View Help

int
id
assign
num
if
(
id
relop
num
)
{
id
assign
num
}
```