

COMP 520 Milestone 2

Rajveer Gandhi, Archit Agnihotri, Dipanjan Dutta

Abstract

This is a design and implementation report of the Milestone 2 of the GoLite project.

1. Design Decisions

Since we got less marks in the Milestone 1 report, we decided to make this report more detailed and descriptive. We decided to include the ideas we missed regarding the scanner, parser and the syntax tree in the first milestone report.

In this milestone Dipanjan created the invalid programs and was responsible for the report. Dipanjan was also responsible for creating the type-checker and the weeder that we could not implement properly in the first milestone. Archit and Raj did the script for symbol table generation and fixed the AST and Pretty printer from Milestone 1.

2. Flex and Bison

In the last report, we mentioned that we chose flex and bison because all of us are comfortable with programming in C. Flex and Bison are frameworks in C that can generate code for C. Our program mainly consists of the following files :

- `src\golite.l` : Scanner
- `src\golite.y` : Parser
- `src\tree.h` : Header file for syntax tree
- `src\tree.c` : Syntax tree generator
- `src\weed.h` : Header file for the weeder

- 21 • **src\weed.c** : Weeder
- 22 • **src\pretty.h** : Header file for the pretty printer
- 23 • **src\pretty.c** : Pretty printer
- 24 • **src\main.c** : Main file which controls execution
- 25 • **src\Makefile** : Script to run flex and bison to create the lexer and
26 parser
- 27 • **build.sh** : Shell script to clear older versions and run the latest version
28 of **Makefile**
- 29 • **run.sh** : Shell script to run **main.c** on a program file with a specified
30 mode.
- 31 • **test.sh** : Shell script to test programs (used of internal testing only)

32 The available modes are :

- 33 • **"scan"** : Scans the program and gives **OK** if it scans successfully,
34 otherwise throws the error.
- 35 • **"tokens"** : Prints the tokens encountered in order inside the program.
- 36 • **"parse"** : Parses the program, creates the syntax tree and prints **OK**
37 if successful, otherwise exits after printing the parsing error.
- 38 • **"pretty"** : Pretty prints the program from traversing the syntax tree.

39 **3. Scanner**

40 The scanner was developed with a very standard setup. We had macros
41 for decimal, octal and hexadecimal representations, and also for escape se-
42 quences for runes and strings. Tokens were captured with the help of regular
43 expressions. Every operator has its individual token. We decided not to
44 group operators in order to keep the code more transparent and interpretable.

45 One challenge that we faced in the scanner was optional semicolons. We
46 had to deal with this by creating a variable called *lastRead*, and updating it
47 whenever a new token is generated (before it is returned). then, when \n or
48 **EOF** or block comment including a \n is encountered we check to see with

49 function *insertSemicolon* if there should have been a semicolon there, and
50 we insert the semicolon in that case.

51 Another challenge we faced was handling block comments. Dipanjan ini-
52 tially added a state machine implementation of the handling of block com-
53 ments. However, it was too abstract and none of us knew how exactly it
54 was handling the block comments. So we later decided to switch to regular
55 expressions.

56 4. Parser

57 Defining the production terms was a bit complex. We had a union struc-
58 ture with all the types of nodes we can have in our syntax tree. We then
59 typed the production nodes according to these node types. Thus the syntax
60 tree was created according to the grammar rules in the parser.

61 Tokens were aliased with the addition of *t* before their natural names.
62 Almost all the tokens were of type string except integer and float literals
63 which had their respective types. We did not have special tokens for data
64 types and boolean *true* and *false* as they were not reserved keywords in
65 GoLite.

66 For the error message generation and display, we used bison's error re-
67 porting function *yyerror*.

68 Associativity and precedence was also defined in the parser. The rule of
69 precedence chosen was (from lower to higher):

- 70 • = (Assignment operator)
- 71 • || (Logical OR operator)
- 72 • && (Logical AND operator)
- 73 • ==, !=, <, >, <=, >=
- 74 • +, -, |, (XOR operator)
- 75 • *, /, %, <<, >>, &, AND-XOR Operator
- 76 • All unary operators

77 All operators are associated left-to-right, except assignment operator,
78 which is associated right-to-left.

79 The production rules are created by a joint agreement between all three
80 of us. One thing Dipanjan missed was that typecasting was actually treated
81 as a function call and not a separate expression type. Raj caught that and
82 changed it.

83 5. Syntax tree

84 The syntax tree that was in the first milestone had a lot of bugs so
85 Archit decided to create the syntax tree altogether. Previously, we had one
86 single convoluted NODE with different structs for different types. This time
87 around, the node is split up into 23 different types of nodes. It might seem
88 too many nodes but it helps bring clarity to the coding structure and helps
89 traversing the syntax tree intuitive. The different types of nodes that are
90 defined are as follows:

- 91 • *PROGRAM* : Program node, root or beginning of the tree.
- 92 • *PACKAGE* : Package declaration.
- 93 • *TOplevelDECL* : Top level declarations. Can be one of the fol-
94 lowing:
 - 95 – *FUNCDCL* : Function declaration
 - 96 – *DCL* : Non-function declaration. They can be :
 - 97 * *VARDCCL* : Top level variable declaration
 - 98 * *TYPEDCL* : Top level type declaration
- 99 • *FUNC_SIGNATURE* : Function signature.
- 100 • *PARAM_LIST* : Parameter list for a function signature.
- 101 • *IDLIST* : List of identifiers for short declarations, assignments and
102 top-level declarations.
- 103 • *TYPE* : Type of a declaration.
- 104 • *STRUCT_TYPE* : Block of the structure.
- 105 • *BLOCK* : Block of code enclosed by braces.

- 106 • *STATEMENTS* : Node acting as a linked list for consecutive state-
107 ments.
- 108 • *STATEMENT* : Program statement.
- 109 • *ELSE_BLOCK* : Special node for *else* part of an *if* block.
- 110 • *SWITCH_CONDITION* : Condition for a switch block.
- 111 • *SWITCH_CASELIST* : Linked-list like node for all the cases and
112 default case for the switch block.
- 113 • *EXPRLIST* : Linked-list like node for expressions (for multiple as-
114 signments in one line).
- 115 • *FOR_CONDITION* : Node for the condition of a for loop.
- 116 • *SIMPLE* : Simple statement.
- 117 • *OTHER_EXPR* : Other statements (including function calls, struct
118 member selector and slice indexing).

119 6. Weeder

120 The weeder is implemented in order to catch any syntactic errors that
121 might not have been directly caught in the parser's grammar. Dipanjan
122 created the weeder on his own. The following checks were implemented in
123 the weeder:

124 6.1. Single default case

125 The language specifies that there can only be one default case inside a
126 single switch block. We used a local variable *hasDefault* that is initailly set to
127 zero. If a default case is encountered, it is set to one. If another default case
128 is found while *hasDefault* is 1, we throw an error. We used a local variable
129 to correclt weed through nested switch cases which are recursive calls. Once
130 the switch block is complete, *hasDefault* is set back to zero.

131 6.2. Break and continue statements

132 According to the language specification, *continue* statements can only
133 appear inside a for loop and *break* statements can only appear inside a for
134 loop or switch block. We had two global static integer variables *insideFor* and
135 *insideSwitch* that stores the level of depth inside a for loop and switch block,
136 respectively (Initial values of both are zero). When a for loop is encountered,
137 *insideFor* is incremented by 1 and control enters the for loop. When control
138 comes back after finishing the for loop, *insideFor* is decremented by 1. The
139 behaviour is similar for *insideSwitch* as well. If a *continue* statement is
140 encountered when *insideFor* is zero, an error is thrown. If a *break* statement
141 is encountered when both *insideFor* and *insideSwitch* are zero, an error is
142 thrown. Making *insideFor* and *insideSwitch* global and integer helps weed
143 nested loops and nested blocks in a program.

144 6.3. Blank identifier

145 According to Go specification, blank identifiers can be used :

- 146 • as an identifier in a declaration
- 147 • as an operand on the left side of an assignment
- 148 • as an identifier on left side of = assignment only
- 149 • as a parameter name in function declaration

150 We used a global boolean variable *isBlankIdValid* that is initially set to false.
151 It is set to true before :

- 152 • the weeding of identifier lists of a declaration
- 153 • *LHS* of an assignment statement
- 154 • *LHS* of a short declaration statement

155 Whenever the recursive weeding of the above are done, the *isBlankIdValid*
156 variable is set to false. If during any expression or statement evaluation, we
157 encounter an "_" we check if *isBlankIdValid* is true or false. If it is false, we
158 throw an error.

159 6.4. Unbalanced assignments

160 If the *LHS* and the *RHS* of an assignment or declaration statement
161 contain unequal number of operands, the compile should quit with appropriate error. This is implemented with two local variables *lhsCount* and
162 *rhsCount*. The *weedIDLIST* and *weedEXPRLIST* functions return the number of operands in the idlist in the *LHS* of an assignment or declaration, or
163 the number of expressions in the *RHS* of the same. The counters are stored in *lhsCount* and *rhsCount* respectively. If they do not match, an error is
164 thrown. Otherwise, the counters are reset to zero for the next statement.
165
166
167

168 6.5. Return statements

169 The language specifies that return statements can only be inside a function body. We use a static global integer variable *insideFunction* that is
170 initially set to zero. When we are weeding the function declaration, we increment the *insideFunction* variable by one and then start weeding of the
171 function block. As soon as weeding of the block is done and control comes back, we decrement *insideFunction* by one. If a return statement is encountered and the value of *insideFunction* is zero, an error is thrown. Use of
172 an integer variable allows for weeding of nested functions (which are not supported in the language yet, but its a nice provision to have).
173
174
175
176
177

178 7. Pretty Printer

179 The pretty printer is almost similar in coding structure to the syntax tree. Besides the fact that we implemented a different node structure for syntax
180 tree, we encountered multiple segmentation faults in Milestone 1 and several other issues. For this reason, we decided to a complete overhaul of the pretty printer. This was done by Archit.
181
182
183

184 Previously, we had a single function split into cases based on the kind of statement/expression inside one function. This time around, we created
185 separate functions corresponding to each type of node in the syntax tree for the pretty printer. The *PROGRAM* is the root node from where the pretty printer starts. Semicolons are printed at the end of every statement.
186
187
188

189 We also have a separate function, *prettyIndent*, that prints 4 spaces per level of indentation level. Indentation level is stored in a variable *g_indent*.
190 Every time we enter a block for pretty printing, the *g_indent* variable increments by one. On exiting the block, *g_indent* is decremented. We decided
191 to print spaces instead of tabs because we have see tabs to create chaotic
192
193

¹⁹⁴ indentations in several editors. Thus to keep a standard across all platforms
¹⁹⁵ for pretty printing, we used spaces.