

CSC4180 Assignment 1 Report

Name: Guangxin Zhao

Student ID: 120090244

Date: Feb 22nd, 2024

Code Structure

```
1 project
2 | -
3 | --- README.md
4 | -
5 | --- testcases
6 | -
7 | --- src
8     | -
9     | --- ir_generator.cpp
10    | --- ir_generator.hpp
11    | --- main.cpp
12    | --- Makefile
13    | --- node.cpp
14    | --- node.hpp
15    | --- parser.y
16    | --- scanner.l
```

How to execute the compiler?

Use commands below to run the compiler:

```
1 cd /path/to/project
2 make all
3
4 120090244@c2d52c9b1339:~/CSC4180-Compiler/Assignment1/src$ ./compiler --help
5 CUHK-SZ CSC4180 Assignment-1: Micro Language Compiler Frontend
6 Usage: Usage: compiler [options] source-program.m
7 Allowed options:
8   -h [ --help ]           Usage: compiler [options]
9                             source-program.m
10  -s [ --scan-only ]       [Default: false] print out token class and
11                             lexeme pairs for each token, no parsing
12                             operations onwards
13  -c [ --cst-only ]        [Default: false] generate concrete syntax
14                             tree only, do not generate AST and LLVM IR
```

```

15  -d [ --dot ] arg (=ast.dot)      [Default: ast.dot] the .dot filename where
16                                     compiler will output the tree
17  -o [ --output ] arg (=program.ll) [Default: program.ll] LLVM IR file compiled
18                                     from source code
19  --source-program arg              source Micro program to compile

```

Sample: test0.m

```

1  120090244@c2d52c9b1339:~/CSC4180-Compiler/Assignment1/src$ ./compiler
  ../testcases/test0.m
2  120090244@c2d52c9b1339:~/CSC4180-Compiler/Assignment1/src$ dot -Tpng ./ast.dot -o
  ./ast.png
3  120090244@c2d52c9b1339:~/CSC4180-Compiler/Assignment1/src$ opt ./program.ll -S --O3 -o
  ./program_optimized.ll
4  120090244@c2d52c9b1339:~/CSC4180-Compiler/Assignment1/src$ llc -march=riscv64
  ./program_optimized.ll -o ./program.s
5  120090244@c2d52c9b1339:~/CSC4180-Compiler/Assignment1/src$ riscv64-unknown-linux-gnu-gcc
  ./program.s -o ./program
6  120090244@c2d52c9b1339:~/CSC4180-Compiler/Assignment1/src$ qemu-riscv64 -L
  /opt/riscv/sysroot ./program
7  30

```

How do I design the Scanner?

The scanner is designed to extract tokens of Micro language according to regular expressions. There are totally 14 tokens in Micro and the Regular Expression rules I designed are addressed below:

```

1  EOLN      "\n"
2  TAB       "\t"
3  COMMENT   "--.*\n"
4  BEGIN_    "begin"
5  END       "end"
6  READ      "read"
7  WRITE     "write"
8  LPAREN    "("
9  RPAREN    ")"
10 SEMICOLON ";"
11 COMMA     ","
12 ASSIGNOP  ":@"
13 PLUSOP    "+"
14 MINUSOP   "-"
15 ID        [a-zA-Z][a-zA-Z0-9_]{0,31}
16 INTLITERAL -?[0-9]+

```

Extracting the tokens using Regular Expressions, we can get tokens stored with type `yylval.strval` or `yylval.intval`, and pass to the parser as `T_TOKENNAME`.

The scanner is also able to print both token class and lexeme (i.e. `<token-class, lexeme>`) for each token with `--scan-only` option. Here is a sample output for the print function on `test0.m`:

```
1 <BEGIN_, begin>
2 <ID, A>
3 <ASSIGNOP, :=>
4 <INTLITERAL, 10>
5 <SEMICOLON, ;>
6 <ID, B>
7 <ASSIGNOP, :=>
8 <ID, A>
9 <PLUSOP, +>
10 <INTLITERAL, 20>
11 <SEMICOLON, ;>
12 <WRITE, write>
13 <LPAREN, (>
14 <ID, B>
15 <RPAREN, )>
16 <SEMICOLON, ;>
17 <END, end>
18 <SCANEOF>
```

How do I design the Parser?

The Parser is designed to receive the tokens extracted from scanner. It also generates a parse tree (or concrete syntax tree), and furthermore, the abstract syntax tree (AST) based on the context-free grammar (CFG).

Here is the CFG of Micro language:

```
1 <start> → <program> SCANEOF
2 <program> → BEGIN <statement list> END
3 <statement list> → <statement> { <statement> }
4 <statement> → ID ASSIGNOP <expression>;
5 <statement> → READ LPAREN <id list> RPAREN;
6 <statement> → WRITE LPAREN<expr list> RPAREN;
7 <id list > → ID { COMMA ID }
8 <expr list > → <expression> { COMMA <expression> }
9 <expression> → <primary> { <add op> <primary> }
10 <primary> → LPAREN <expression> RPAREN
11 <primary> → ID
12 <primary> → INTLITERAL
13 <add op> → PLUSOP
14 <add op> → MINUSOP
```

In `parser.y`, I designed three `yyval` data types:

```

1  %union {
2      /* Data */
3      int intval;
4      const char* strval;
5      struct Node* nodeval;
6  };

```

Then, assign these data types to terminal symbols and non-terminal symbols.

Terminal Symbols:

```

1  %token <intval> T_INTLITERAL
2  %token <strval> T_BEGIN_ T_END T_READ T_WRITE T_LPAREN T_RPAREN T_SEMICOLON
3  %token <strval> T_COMMA T_ASSIGNOP T_PLUSOP T_MINUSOP T_SCANEOF T_ID

```

Non-Terminal Symbols:

```

1  %type <nodeval> program statement_list statement id_list expr_list expression primary

```

How do I design the Intermediate Code Generator?

The intermediate representation (IR) of the compiler is the LLVM IR, and the IR generator converts the AST given from the parser to LLVM IR (.ll) file. The `ir_generator.cpp` file has the following structure:

Export AST to LLVM IR

The `export_ast_to_llvm_ir()` function generates the LLVM IR header and main function structure. And it calls the `gen_llvm_ir()` function.

Generate LLVM IR

The generate LLVM IR sequence recursively calls the main routine or the sub-routines to generate the full LLVM IR.

- Function `gen_llvm_ir()`

This is the main routine in the generation part, it calls sub-routines such as `gen_assignop_llvm_ir()`, `gen_read_llvm_ir()`, `gen_write_llvm_ir()`, and also, the main routine `gen_llvm_ir()` to traverse the whole AST tree.

- Function `gen_assignop_llvm_ir()`

This sub-routine handles the following LLVM IR instructions:

```

1  %<variable> = alloca i32
2  store i32 <rvalue>, i32* %<variable>

```

- Function `gen_read_llvm_ir()`

This sub-routine handles the following LLVM IR instructions:

```
1 | %<variable> = alloca i32
2 | %_scanf_format_1 = alloca [# x i8]
3 | store [# x i8] c"%d ... %d\00", [# x i8]* %_scanf_format_1
4 | %_scanf_str_1 = getelementptr [# x i8], [# x i8]* %_scanf_format_1, i32 0, i32 0
5 | call i32 (i8*, ...) @scanf(i8* %_scanf_str_1, i32* %<variable>)
```

- Function `gen_write_llvm_ir()`

This sub-routine handles the following LLVM IR instructions:

```
1 | %_printf_format_1 = alloca [# x i8]
2 | store [# x i8] c"%d ... %d\0A\00", [# x i8]* %_printf_format_1
3 | %_printf_str_1 = getelementptr [# x i8], [# x i8]* %_printf_format_1, i32 0, i32 0
4 | call i32 (i8*, ...) @printf(i8* %_printf_str_1, i32 <rvalue>)
```

- Function `gen_operation_llvm_ir()`

This sub-routine handles the following LLVM IR instructions:

```
1 | %_tmp_1 = load i32, i32* %<variable>
2 | %_tmp_2 = sub i32 <expression>, %_tmp_1
3 | %_tmp_3 = add i32 %_tmp_1, %_tmp_2
```

- Function `gen_expression_llvm_ir()`

This sub-routine handles the following LLVM IR instruction:

```
1 | %_tmp_1 = load i32, i32* %<variable>
```