

# CSE 340 Spring 2017 – Project 4

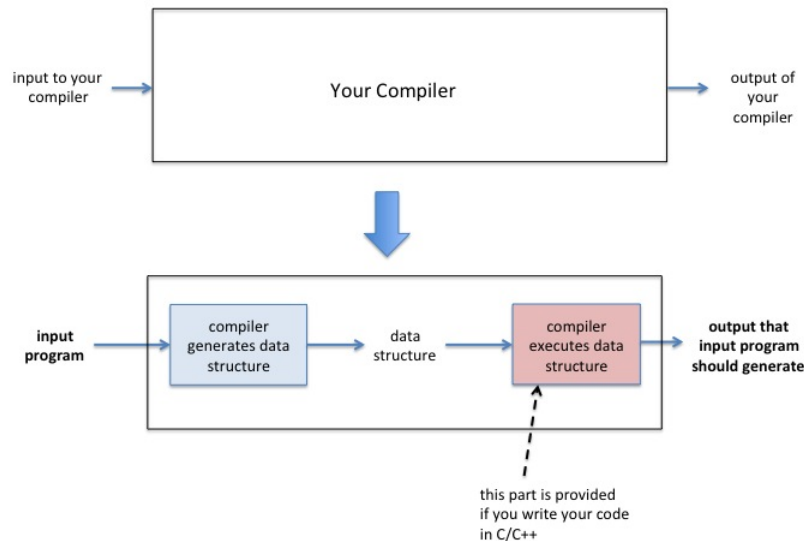
Due on **April 28, 2017** by **11:59 pm**

## Abstract

The goal of this project is to give you some hands-on experience with implementing a small compiler. You will write a compiler for a simple language. You will not be generating low level code. Instead, you will generate an intermediate representation (a data structure that represents the program). The execution of the program will be done after compilation by *interpreting* the generated intermediate representation.

## 1 Introduction

You will write a small compiler that will read an input program and represent it in an internal data structure. The data structure will contain representation of instructions to be executed as well as a part that represents the memory of the program (locations for variables). Then your compiler will execute the data structure (interpret it). This means that the program will traverse the data structure and at every node it visits, it will execute the node by changing the content of memory locations and deciding what is the next instruction to execute (program counter). The output of your compiler is the output that the input program should produce. These steps are illustrated in the following figure



The remainder of this document is organized as follows:

1. **Grammar** Defines the programming language syntax including grammar.

2. **Execution Semantics** Describe statement semantics for `if`, `while`, `switch` and `print` statements.
3. **How to generate the intermediate representation** Explains step by step how to generate the intermediate representation (data structure). You should read this sequentially and not skip around.
4. **Executing the intermediate representation** Basically, you have two options. If you are using C/C++, you should use the code we provide to execute the intermediate representation. If you are using Java, it describes the strict rules to follow in executing the intermediate representation. Those rules will be enforced.
5. **Requirements** Lists the programming languages you are allowed to use in your solution (C/C++ or Java) and other requirements.
6. **Grading** Describes the grading scheme.
7. **Bonus Project** Describes the requirements for the bonus project.

## 2 Grammar

The grammar for this project is a simplified form of the grammar from the previous project, but there are a couple extensions.

<i>program</i>	→	<i>var_section body</i>
<i>var_section</i>	→	<i>id_list</i> SEMICOLON
<i>id_list</i>	→	ID COMMA <i>id_list</i>   ID
<i>body</i>	→	LBRACE <i>stmt_list</i> RBRACE
<i>stmt_list</i>	→	<i>stmt stmt_list</i>   <i>stmt</i>
<i>stmt</i>	→	<i>assign_stmt</i>   <i>print_stmt</i>   <i>while_stmt</i>   <i>if_stmt</i>   <i>switch_stmt</i>   <i>for_stmt</i>
<i>assign_stmt</i>	→	ID EQUAL <i>primary</i> SEMICOLON
<i>assign_stmt</i>	→	ID EQUAL <i>expr</i> SEMICOLON
<i>expr</i>	→	<i>primary op primary</i>
<i>primary</i>	→	ID   NUM
<i>op</i>	→	PLUS   MINUS   MULT   DIV
<i>print_stmt</i>	→	<b>print</b> ID SEMICOLON
<i>while_stmt</i>	→	WHILE <i>condition body</i>
<i>if_stmt</i>	→	IF <i>condition body</i>
<i>condition</i>	→	<i>primary relop primary</i>
<i>relop</i>	→	GREATER   LESS   NOTEQUAL
<i>switch_stmt</i>	→	SWITCH ID LBRACE <i>case_list</i> RBRACE
<i>switch_stmt</i>	→	SWITCH ID LBRACE <i>case_list default_case</i> RBRACE
<i>for_stmt</i>	→	FOR LPAREN <i>assign_stmt condition</i> SEMICOLON <i>assign_stmt</i> RPAREN <i>body</i>
<i>case_list</i>	→	<i>case case_list</i>   <i>case</i>
<i>case</i>	→	CASE NUM COLON <i>body</i>
<i>default_case</i>	→	DEFAULT COLON <i>body</i>

### Some highlights of the grammar:

1. Expressions are greatly simplified and are not recursive.
2. There is no type declaration section.
3. Division is integer division and the result of the division of two integers is an integer.
4. *if* statement is introduced. Note that *if\_stmt* does not have *else*.
5. *for* statement is introduced. Note that *for* has a very general syntax similar to that of the C language for loop.
6. A *print* statement is introduced. Note that the **print** keyword is in lower case, but other keywords are all upper-case letters.
7. There is no variable declaration list. There is only one *id\_list* in the global scope and that contains all the variables.
8. There is no type specified for variables. All variables are **int** by default.

## 3 Execution Semantics

All statements in a statement list are executed sequentially according to the order in which they appear. Exception is made for body of *if\_stmt*, *while\_stmt* and *switch\_stmt* as explained below. In what follows, I will assume that all values are stored in locations even constants. This assumption is used by the execution procedure that we provide. This is not a restrictive assumption. For variables, you will have locations associated with them. For constants, you can reserve a location in which you store the constant (this is like having an unnamed immutable variable).

### 3.1 Assignment Statement

To execute an assignment statement, the expression on the lefthand side of the equal sign is evaluated and the result is stored in the location associated with the righthand side of the expression.

### 3.2 Expression

To evaluate an expression, the values in the locations associated with the two operands are obtained and the expression operator is applied to these values resulting in a value for the expression.

### 3.3 Boolean Condition

A boolean condition takes two operands as parameters and returns a boolean value. It is used to control the execution of *while* and *if* statements. To evaluate a condition, the values in the locations associated with the operands are obtained and the relational operator is applied to these values resulting in a true or false value. For example, if the values of the two operands *a* and *b* are 3 and 4 respectively, *a < b* evaluates to **true**.

### 3.4 *If* statement

*if\_stmt* has the standard semantics:

1. The condition is evaluated.
2. If the condition evaluates to **true**, the body of the *if\_stmt* is executed, then the next statement following the *if* in the *stmt\_list* is executed.
3. If the condition evaluates to **false**, the statement following the *if* in the *stmt\_list* is executed

These semantics apply recursively to nested *if\_stmt*.

### 3.5 *While* statement

*while\_stmt* has the standard semantics:

1. The condition is evaluated.
2. If the condition evaluates to **true**, the body of the *while\_stmt* is executed, then the condition is evaluated again and the process repeats.
3. If the condition evaluates to **false**, the statement following the *while\_stmt* in the *stmt\_list* is executed.

These semantics apply recursively to nested *while\_stmt*. The code block:

```
WHILE condition
{
    stmt_list
}
```

is equivalent to:

```
label: IF condition
{
    stmt_list
    goto label
}
```

Note that **goto** statements do not appear in the input program, but our intermediate representation includes **GotoStatement** which is used in conjunction with **IfStatement** to represent *while* and *switch* statements.

### 3.6 *For* statement

The *for\_stmt* is very similar to the C language for statement. The semantics are defined by giving an equivalent construct.

```
FOR ( assign_stmt_1 condition ; assign_stmt_2 )  
{  
    stmt_list  
}
```

is equivalent to:

```
assign_stmt_1  
WHILE condition  
{  
    stmt_list  
    assign_stmt_2  
}
```

For example, the following snippet of code:

```
FOR ( a = 0; a < 10; a = a + 1; )  
{  
    print a;  
}
```

is equivalent to:

```
a = 0;  
WHILE a < 10  
{  
    print a;  
    a = a + 1;  
}
```

### 3.7 *Switch* statement

*switch\_stmt* has the following semantics<sup>1</sup>:

1. The value of the switch variable is checked against each case number in order.
2. If the value matches the number, the body of the case is executed, then the statement following the *switch\_stmt* in the *stmt\_list* is executed.

---

<sup>1</sup>The switch statement in the C language has different syntax and semantics. It is also more dangerous!

3. If the value does not match the number, next case is evaluated.
4. If a default case is provided and the value does not match any of the case numbers, then the body of the default case is executed and then the statement following the *switch\_stmt* in the *stmt\_list* is executed.
5. If there is no default case and the value does not match any of the case numbers, then the statement following the *switch\_stmt* in the *stmt\_list* is executed.

These semantics apply recursively to nested *switch\_stmt*. The code block:

```
SWITCH var {
    CASE  $n_1$  : { stmt_list_1 }
    ...
    CASE  $n_k$  : { stmt_list_k }
}
```

is equivalent to:

```
IF var ==  $n_1$  {
    stmt_list_1
    goto label
}
...
IF var ==  $n_k$  {
    stmt_list_k
    goto label
}
label:
```

And for switch statements with default case, the code block:

```
SWITCH var {
    CASE  $n_1$  : { stmt_list_1 }
    ...
    CASE  $n_k$  : { stmt_list_k }
    DEFAULT : { stmt_list_default }
}
```

is equivalent to:

```

    IF var == n1 {
        stmt_list_1
        goto label
    }
    ...
    IF var == nk {
        stmt_list_k
        goto label
    }
    stmt_list_default
label:

```

### 3.8 *Print* statement

The statement

```
print a;
```

prints the value of variable *a* at the time of the execution of the *print* statement.

## 4 How to generate the code

The intermediate code will be a data structure (a graph) that is easy to interpret and execute. I will start by describing how this graph looks for simple assignments then I will explain how to deal with *while* statements.

Note that in the explanation below I start with incomplete data structures then I explain what is missing and make them more complete. You should read the whole explanation.

### 4.1 Handling simple assignments

A simple assignment is fully determined by: the operator (if any), the id on the left-hand side, and the operand(s). A simple assignment can be represented as a node:

```

struct AssignmentStatement {
    struct ValueNode* left_hand_side;
    struct ValueNode* operand1;
    struct ValueNode* operand2;
    ArithmeticOperatorType op; // operator
}

```

The operator of an assignment statement without an expression on the righthand side is set to `OPERATOR_NONE` and there is only one operand. To execute an assignment, you need the values of the operand(s), apply the operator, if any, to the operands and assign the resulting value of the right-hand side to the `left_hand_side`. For literals (`NUM`), the value is the value of the number. For variables, the value is the last value stored in the location associated with the variable. **Initially, all variables are initialized to 0.**

Multiple assignments are executed one after another. So, we need to allow multiple assignment nodes to be linked to each other. This can be achieved as follows:

```
struct AssignmentStatement {
    struct ValueNode* left_hand_side;
    struct ValueNode* operand1;
    struct ValueNode* operand2;
    ArithmeticOperatorType op; // operator
    struct AssignmentStatement* next;
}
```

This structure only accepts `ValueNode` as operands. To handle literal constants (`NUM`), you need to create `ValueNode` or them and initialize the value stored in these value nodes to the constant value. This initialization, as well as the initialization of values in locations associated with variable is done while parsing.

This will now allow us to execute a sequence of assignment statements represented in a linked-list: we start with the head of the list, then we execute every assignment in the list one after the other.

**Begin Note** It is important to distinguish between compile-time initialization and runtime execution. For example, consider the program

```
a, b;
{
    a = 3;
    b = 5;
}
```

The intermediate representation for this program will have have two assignment statements: one to copy the value in the location that contains the value 3 to the location associated with a and one to copy the value in the location that contains the value 5 to the location associated with b. The values 3 and 5 will not be copied to the locations of a and b at compile-time. **End Note**

This is simple enough, but does not help with executing other kinds of statements. We consider them one at a time.

## 4.2 Handling *print* statements

The *print* statement is straightforward. It can be represented as

```
struct PrintStatement
{
    struct ValueNode* id;
}
```

Now, we ask: how can we execute a sequence of statements that are either assignments or print statements (or other types of statements)? We need to put both kinds of statements in a list and not just the assignment statements as we did above. So, we introduce a new kind of node: a statement node. The statement node has a field that indicates which type of statement it is. It also has fields to accommodate the remaining types of statements. It looks like this:



```

struct StatementNode {
    StatementType type; // NOOP_STMT, GOTO_STMT, ASSIGN_STMT, IF_STMT, PRINT_STMT

    union {
        struct AssignmentStatement* assign_stmt;
        struct PrintStatement* print_stmt;
        struct IfStatement* if_stmt;
        struct GotoStatement* goto_stmt;
    };
    struct StatementNode* next;
}

```

This way we can go through a list of statements and execute one after the other. To execute a particular node, we check its `type`. If it is `PRINT_STMT`, we execute the `print_stmt` field, if it is `ASSIGN_STMT`, we execute the `assign_stmt` field and so on. With this modification, we do not need a `next` field in the `AssignmentStatement` structure (as we explained above), instead, we put the `next` field in the statement node.

This is all fine, but we do not yet know how to generate the list to execute later. The idea is to have the functions that parse non-terminals return the code for the non-terminals. For example for a statement list, we have the following pseudocode (missing many checks):

```

struct StatementNode* parse_stmt_list()
{
    struct StatementNode* st; // statement
    struct StatementNode* stl; // statement list

    st = parse_stmt();
    if (nextToken == start of a statement list)
    {
        stl = parse_stmt_list();
        append stl to st; // this is pseudocode
        return st;
    }
    else
    {
        ungetToken();
        return st;
    }
}

```

And to parse *body* we have the following pseudocode:

```

struct StatementNode* parse_body()
{
    struct StatementNode* stl;

    match LBRACE
    stl = parse_stmt_list();
    match RBRACE

    return stl;
}

```

### 4.3 Handling *if* and *while* statements

More complications occur with *if* and *while* statements. The structure for an *if* statement can be as follows:

```
struct IfStatement {
    ConditionalOperatorType condition_op;
    struct ValueNode* condition_operand1;
    struct ValueNode* condition_operand2;

    struct StatementNode* true_branch;
    struct StatementNode* false_branch;
}
```

The `condition_op`, `condition_operand1` and `condition_operand2` fields are the operator and operands of the condition of the *if* statement. To generate the node for an *if* statement, we need to put together the *condition*, and *stmt\_list* that are generated in the parsing of the *if* statement.

The `true_branch` and `false_branch` fields are crucial to the execution of the *if* statement. If the condition evaluates to true then the statement specified in `true_branch` is executed otherwise the one specified in `false_branch` is executed. We need one more type of node to allow loop back for *while* statements. This is a *GotoStatement*.

```
struct GotoStatement {
    struct StatementNode* target;
}
```

To generate code for *while* and *if* statements, we need to put a few things together. The outline given above for *stmt\_list* needs to be modified as follows (this is missing details and shows only the main steps).

```
struct StatementNode* parse_stmt()
{
    ...

    create statement node st
    if next token is IF
    {
        st->type = IF_STMT;
        create if-node;                                // note that if-node is pseudocode and is not
                                                         // a valid identifier in C, C++ or Java

        st->if_stmt = if-node;

        parse the condition and set if-node->condition_op, if-node->condition_operand1 and if-node->condition_operand2

        if-node->true_branch = parse_body();             // parse_body returns a pointer to a list of statements

        create no-op node                               // this is a node that does not result
                                                         // in any action being taken

        append no-op node to the body of the if         // this requires a loop to get to the end of
                                                         // if-node->true_branch by following the next field
                                                         // you know you reached the end when next is NULL
                                                         // it is very important that you always appropriately
                                                         // initialize fields of any data structures
                                                         // do not use uninitialized pointers

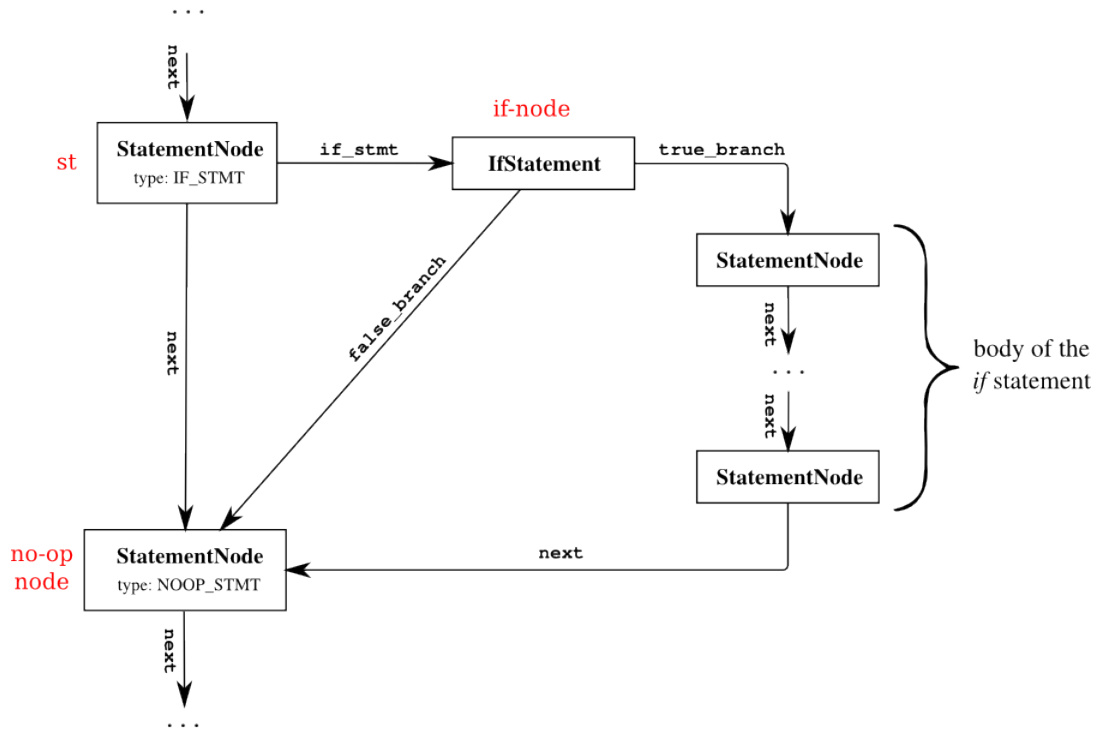
        set if-node->>false_branch to point to no-op node
        set st->next to point to no-op node
    }
}
```

```

...
} else ...
}

```

The following diagram shows the desired structure for the *if* statement:



The *stmt\_list* code should be modified because of the extra no-op node:

```

struct StatementNode* parse_stmt_list()
{
    struct StatementNode* st; // statement
    struct StatementNode* stl; // statement list

    st = parse_stmt();
    if (nextToken == start of a statement list)
    {
        stl = parse_stmt_list();

        if st->type == IF_STMT
        {
            append stl to the no-op node that follows st

            //      st
            //      |
            //      v
            //      no-op
            //      |
            //      v
            //      stl
        }
    }
    else
    {

```

```

        append st1 to st;

        //      st
        //      |
        //      V
        //      st1
    }
    return st;
}
else
{
    ungetToken();
    return st;
}
}

```

Handling *while* statement is similar. Here is the outline for parsing a *while* statement and creating the data structure for it:

```

...

create statement node st
if next token is WHILE
{
    st->type = IF_STMT;                // handling WHILE using if and goto nodes
    create if-node                    // if-node is not a valid identifier see
                                     // corresponding comment above
    st->if_stmt = if-node

    parse the condition and set if-node->condition_op, if-node->condition_operand1 and if-node->condition_operand2

    if-node->>true_branch = parse_body();

    create a new statement node gt    // This is of type StatementNode
    gt->type = GOTO_STMT;
    create goto-node                  // This is of type GotoStatement
    gt->goto_stmt = goto-node;
    goto-node->target = st;            // to jump to the if statement after
                                     // executing the body

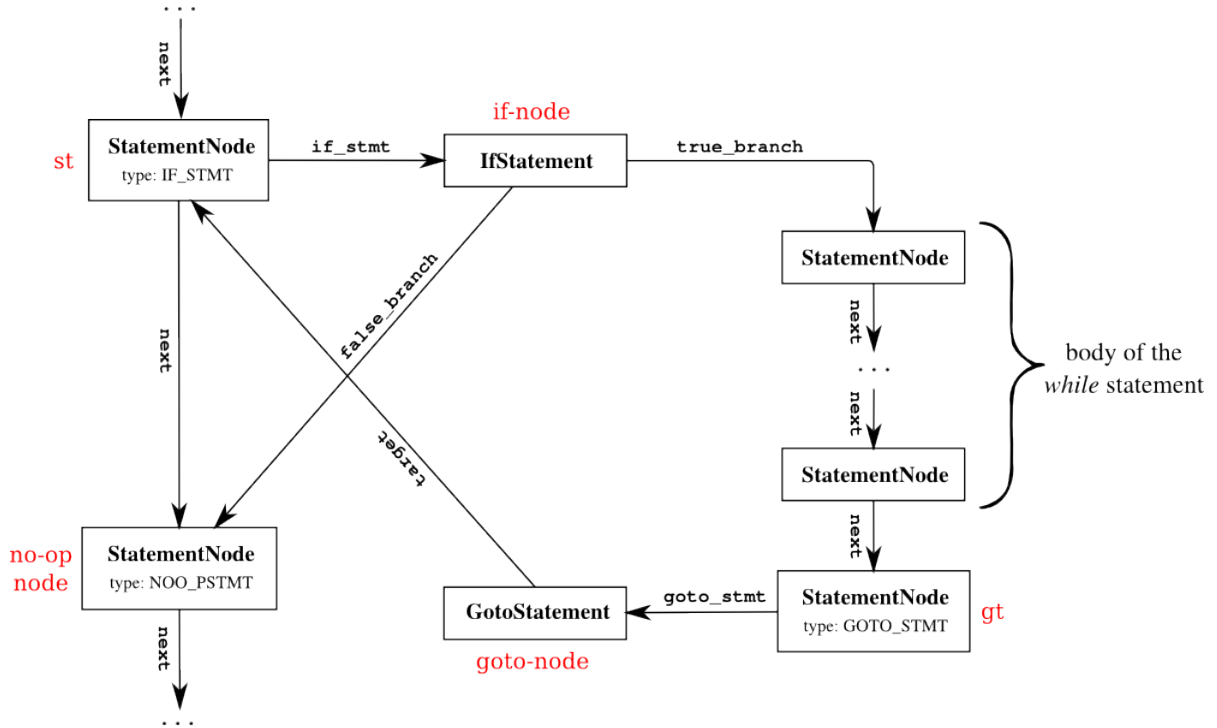
    append gt to the body of the while // append gt to the body of the while
                                     // this requires a loop. check the comment
                                     // for the if above.

    create no-op node
    set if-node->>false_branch to point to no-op node
    set st->next to point to no-op node
}

...

```

The following diagram shows the desired structure for the *while* statement:



#### 4.4 Handling *switch* statement

You can handle the *switch* statement similarly. Use a combination of **IfStatement** and **GotoStatement** to support the semantics of the *switch* statement. See section 3.7 for the semantics of the *switch* statement.

### 5 Executing the intermediate representation

After the graph data structure is built, it needs to be executed. Execution starts with the first node in the list. Depending on the type of the node, the next node to execute is determined. The general form for execution is illustrated in the following pseudo-code.

```
pc = first node
while (pc != NULL)
{
    switch (pc->type)
    {
        case ASSIGN_STMT: // code to execute pc->assign_stmt ...
            pc = pc->next

        case IF_STMT:      // code to evaluate condition ...
                           // depending on the result
                           // pc = pc->if_stmt->>true_branch
                           // or
                           // pc = pc->if_stmt->>false_branch
```

```

    case NOOP_STMT:    pc = pc->next

    case GOTO_STMT:    pc = pc->goto_stmt->target

    case PRINT_STMT:  // code to execute pc->print_stmt ...
                      pc = pc->next
}
}

```

**JAVA Implementation: Executing the graph should be done non-recursively and without any function calls. Even helper functions are not allowed for the execution of the graph. This is a requirement that will be checked by inspecting your code. Little credit will be assigned if this requirement is not met.**

**C/C++ implementations** If you are developing in C/C++, we have provided you with the data structures and the code to execute the graph and **you must use it**. There are two files `compiler.h` and `compiler.cc`. You need to write your code in separate file(s) and `#include compiler.h`. The entry point of your code is a function declared in `compiler.h`:

```
struct StatementNode* parse_generate_intermediate_representation();
```

You need to implement this function.

The `main()` function is provided in `compiler.cc`:

```

int main()
{
    struct StatementNode * program;
    program = parse_generate_intermediate_representation();
    execute_program(program);
    return 0;
}

```

It calls the function that you will implement which is supposed to parse the program and generate the intermediate representation, then it calls the `execute_program` function to execute the program. You should not modify any of the given code. In fact if you write your program in C/C++, you should only submit the file(s) that contain your own code and we will add the given part and compile the code before testing. If you write your program in Java, you should strictly follow the guidelines for executing the intermediate representation.

## 6 Requirements

1. Write a compiler that generates intermediate representation for the code and write an interpreter to execute the intermediate representation. The interpreter is provided for C/C++ implementations.
2. **Language:** You can use Java, C, or C++ for this assignment.
3. **Any language other than Java, C or C++ is not allowed for this project.**

4. If you use C/C++ for this project, you should use the provided code and only implement the required functions.
5. If you use Java, you will need to write everything yourself. The requirements on how to execute the intermediate representation will be checked manually when grading.
6. **Platform:** As in previous projects, the reference platform is CentOS 7
7. **You can assume that there are no syntax or semantic errors in the input program.**

## 7 Grading

The test cases provided with the project, do not contain any test case for *switch* or *for* statements. However, test cases containing these statements will be included on the submission site. Make sure you test your code extensively with input programs that contain *switch* statements and *for* statements.

The test cases (there will be multiple test cases in each category, each with equal weight) will be broken down in the following way (out of 100 points):

- Assignment statements: 20
- If statements: 25
- While statements: 25
- Switch and For statements: 20
- All statements: 10

Note that all test cases depend on successful implementation of print statements (otherwise your program cannot generate correct output), and if, while, for, and switch statement test cases all depend on assignment statements.

## 8 Bonus Project

### 8.1 Bonus Project Options

You have three options for the bonus project:

1. Resubmit project 2. For this option you are given another chance to submit project 2. The grade you obtain will be reduced by 30% and replaces the grade you already obtained on project 2. So, if your grade for project 2 was 20 and your grade for the replacement is 90, the grade for project 2 will be changed from 20 to 63 ( $63 = 90$  reduced by 30%)
2. Resubmit project 3. For this option you are given another chance to submit project 3. The grade you obtain will be reduced by 30% and replaces the grade you already obtained on project 3. So, if your grade for project 3 was 20 and your grade for the replacement is 100, the grade for project 3 will be changed from 20 to 70 ( $70 = 100$  reduced by 30%)
3. Do a new project (advanced compiler) that replaces the lowest grade between project 2 and project 3. The grade you obtain under this option will replace the lower of the two grades that you obtained on project 2 or 3 (if the grade you obtain on the bonus is lower than what you already got on projects 2 and 3, no replacement is made).

If you make multiple submissions, only the last submission will count.

### 8.2 Advanced Compiler (option 3)

Support the following grammar:

<i>program</i>	→	<i>var_section body</i>
<i>var_section</i>	→	VAR <i>int_var_decl array_var_decl</i>
<i>int_var_decl</i>	→	<i>id_list</i> SEMICOLON
<i>array_var_decl</i>	→	<i>id_list</i> COLON ARRAY LBRAC NUM RBRAC SEMICOLON
<i>id_list</i>	→	ID COMMA <i>id_list</i>   ID
<i>body</i>	→	LBRACE <i>stmt_list</i> RBRACE
<i>stmt_list</i>	→	<i>stmt stmt_list</i>   <i>stmt</i>
<i>stmt</i>	→	<i>assign_stmt</i>   <i>print_stmt</i>   <i>while_stmt</i>   <i>if_stmt</i>   <i>switch_stmt</i>
<i>assign_stmt</i>	→	<i>var_access</i> EQUAL <i>expr</i> SEMICOLON
<i>var_access</i>	→	ID   ID LBRAC <i>expr</i> RBRAC
<i>expr</i>	→	<i>term</i> PLUS <i>expr</i>
<i>expr</i>	→	<i>term</i>
<i>term</i>	→	<i>factor</i> MULT <i>term</i>
<i>term</i>	→	<i>factor</i>
<i>factor</i>	→	LPAREN <i>expr</i> RPAREN
<i>factor</i>	→	NUM
<i>factor</i>	→	<i>var_access</i>
<i>print_stmt</i>	→	<b>print</b> <i>var_access</i> SEMICOLON
<i>while_stmt</i>	→	WHILE <i>condition body</i>



<i>if_stmt</i>	→	IF <i>condition</i> <i>body</i>
<i>condition</i>	→	<i>expr</i> <i>relop</i> <i>expr</i>
<i>relop</i>	→	GREATER   LESS   NOTEQUAL
<i>switch_stmt</i>	→	SWITCH <i>var_access</i> LBRACE <i>case_list</i> RBRACE
<i>switch_stmt</i>	→	SWITCH <i>var_access</i> LBRACE <i>case_list</i> <i>default_case</i> RBRACE
<i>case_list</i>	→	<i>case</i> <i>case_list</i>   <i>case</i>
<i>case</i>	→	CASE NUM COLON <i>body</i>
<i>default_case</i>	→	DEFAULT COLON <i>body</i>

Note that LBRAC is "[" and LBRACE is "{". The former is used for arrays and the latter is used for body. Assume that all arrays are integer arrays and are indexed from 0 to *size* - 1, where *size* is the size of the array specified in the *var\_section* after the ARRAY keyword and between "[" and "]".

The data structures and code that we have provided for the regular assignment will not be enough for the bonus, you will need to modify those to support arrays. Submit *all* code files for the bonus project (including the modified `compiler.h` and `compiler.cc`).

**All restrictions imposed on the execution of the intermediate representation for the regular project apply to the bonus project as well. You are not allowed to call any functions while executing the intermediate representation. You are not allowed to execute the program recursively.**