

Projects of Compiler Principles

CS215

Yaoming ZHU

5140219355

Date Performed: January 18, 2017
Instructor: Professor Li Jiang

1 Objective

In this project, it is required to design and implement a simplified compiler, for a given programming language, namely smallC, which can be regarded as a subset and the core part of C programming language. The compiler should translate smallC source code to MIPS assembly codes.

1.1 Project Logic

Lexical Analyzer The *Flex* tool is used to help design a lexical analyzer generator.

Syntax Analyzer The *Bison*(A.K.A *Yacc*) tool is used to help design an LALR parser generator.

Abstract Syntax Tree After syntax analysis, the AST tree can be generated. It helps to test the output of former process.

Semantic Analyzer Possible semantic errors is checked in this part.

Intermediate Representation Code Translation The AST tree is translated into IR code in this part.

Optimization Some optimizations are applied to IR code in this part, i.e. dead code elimination.

Machine-code Generation The IR code is interpreted into MIPS code in this part.

Project Test Test the project with some smallC code and SPIM.

1.2 Project Environment

The project guidance says "In this project, Linux environment (Ubuntu) is required". Fortunately, with the help of *bash on ubuntu on windows*, this project can be implemented under Windows system.

Furthermore, *Win flex-bison* provided with all flex and bison function to the Windows platform, which means this project can be implemented even under Windows IDE! (And I actually finished most parts of the project with Visual Studio 2015). With IDE support, the debugging process becomes much more efficient and precise.

For setup details, please refer to

<https://sourceforge.net/p/winflexbison/wiki/Visual%20Studio%20custom%20build%20rules/>

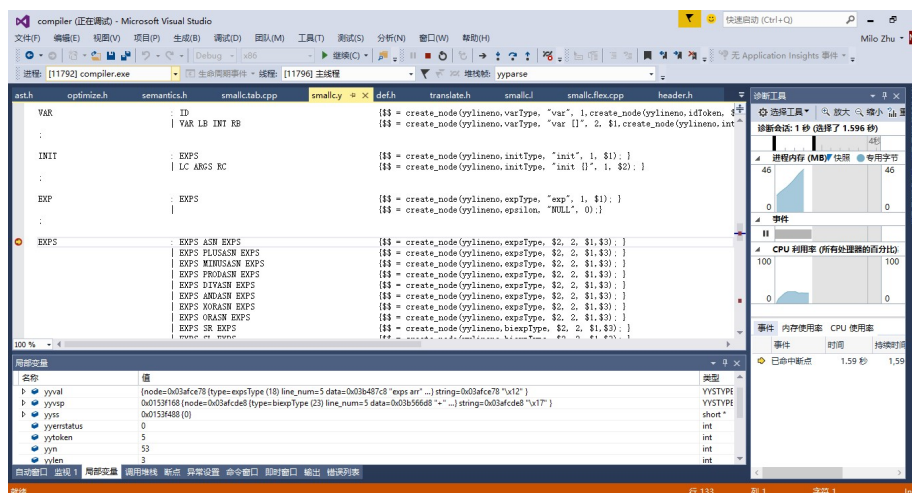


Figure 1: The Project on Visual Studio 2015

Besides the FTP-uploaded one, I have also uploaded a Visual Studio project, they only small difference they have is that, for the output short name of *Win flex-bison* is different from Ubuntu one, thus it leads to different head file reference. You can access Visual Studio version from the link in readMe file.

2 Lexical Analyzer

With the help of *Flex*, this part is fairly simple. This part is done in *smallC.l*.

2.1 Keywords

Since the keywords are fixed, the lexical analyzer will return corresponding token and string value once it find the keywords. *yylval* is provided by *flex* to help to capture the string content of our token.

```
1 int          {yylval.string = strdup(yytext); return TYPE;}
2 struct      {yylval.string = strdup(yytext); return STRUCT;}
3 return      {yylval.string = strdup(yytext); return RETURN;}
4 if          {yylval.string = strdup(yytext); return IF;}
5 else       {yylval.string = strdup(yytext); return ELSE;}
6 break      {yylval.string = strdup(yytext); return BREAK;}
7 continue   {yylval.string = strdup(yytext); return CONT;}
8 for        {yylval.string = strdup(yytext); return FOR;}
9 read       {yylval.string = strdup(yytext); return READ;}
10 write      {yylval.string = strdup(yytext); return WRITE;}
```

Note. According to the guidance, the 'read' and 'write' are not the smallC keywords. However, they are used to test input and output function the compiler, thus they are regarded as keywords in this project.

2.2 Numbers and Identifiers

The token for numbers and identifiers are a little more tricky, regular expression is needed.

```
1 digitnon0 [1-9]
2 hexdigit  [0-9A-Fa-f]
3 letter    [A-Za-z]
4 digit     ({digitnon0}|0)
5
6 identifier ((-|{letter})(-|{letter}|{digit}))*
7 dec_int    (({digitnon0})+({digit}*))
8 oct_int    (0({digit}*))
9 hex_int    (0(x|X)({digit}*))
10 %%
11 ... // keywords
12 {dec_int}  {yylval.string = strdup(yytext); return INT; }
13 {oct_int}  {yylval.string = strdup(yytext); return INT; }
14 {hex_int}  {yylval.string = strdup(yytext); return INT; }
15 {identifier} {yylval.string = strdup(yytext); return ID;}
```

2.3 Blank Characters and Comments

Like most C-like languages, smallC will just ignore the blank characters, so do the lexical analyzer. The only exception is the newline characters: the lexical analyzer will use it to get current line number. A variable *yylineno* is introduced to keep record of current line number.

```
1 extern int yylineno;
2
3 %%
4 [\t\r\f]      {;}
```

```

5 | [\n]                {yylineno++;}
6 | [ ]                 {;}

```

As for comments, two different cases shall be treated, the in-line comment(//) and the out-of-line comment(/* */). The former one will end once the newline character is encountered, the later one will end only if "*/" is encountered.

```

1
2 void commentInLine();
3 void commentOutLine();
4
5 %%
6 "//"      {commentInLine();}
7 "/*"      {commentOutLine();}
8 %%
9
10 void commentInLine(void){
11     char c;
12     while((c = input())!='\n');
13     yylineno++;
14 }
15
16 void commentOutLine(void){
17     char c1, c2;
18     c2 = input();
19     do{
20         c1 = c2;
21         c2 = input();
22         if(c1 == '\n'){
23             yylineno++;
24         }
25     }while(!(c1 == '*' && c2 == '/'));
26 }

```

2.4 Test

Unless otherwise specified, the following smallC code will be used as test code.

```

1 int main()
2 {
3     int x;
4     read(x);
5     x = x + 10;
6     write(x);
7     return 0;
8 }

```

test.c

After *flex* the .l file and compile, we can get the result:
The lexical analyzer works well.

3 Syntax Analyzer

The design of syntax analyzer is simple but tough.

```

F:\complier\project\cs308\Proj
^^ bash
acer@ACER:/mnt/f/complier/project/cs308/Proj$ ./lexscanner <test0
int is token 3
main is token 51
( is token 4
) is token 5
{ is token 8
int is token 3
x is token 51
; is token 1
read is token 52
( is token 4
x is token 51
) is token 5
; is token 1
x is token 51
x is token 51
= is token 40
x is token 51
+ is token 26
10 is token 50
; is token 1
write is token 53
( is token 4
x is token 51
) is token 5
; is token 1
return is token 11
0 is token 50
; is token 1
} is token 9
lex finished!
acer@ACER:/mnt/f/complier/project/cs308/Proj$

```

Figure 2: The Output of Lexical Analyzer

This part is done in *smallC.y*. *Bison* is used in this part.

First and foremost, a data structure is required to represent abstract syntax tree(AST). Thus, a struct is constructed as AST node, which shall have following functions:

- a. Hold the AST node type
- b. Hold the pointers pointing to its sub AST-node.

Thus, it is designed as:

```

1 typedef struct TreeNode {
2     nodeTypeEnum type;
3     int lineNum;
4     char* data;
5     int size, capacity;
6     struct TreeNode** children;
7 } TreeNode;

```

size and *capacity* hold values for number of this kind of AST node and max number of AST can have respectively. *lineNum* is the corresponding line number of the AST node. *data* store some informations. For terminals(i.e. identifier), it contains string content; for non-terminals, it store the derivation case, for most non-terminals have more than one production. *type* is an enum type to hold the AST node type, which is defined as

```

1 typedef enum{
2   programType,
3   extdefsType,
4   extdefType,
5   // ..... other non-terminals
6   unaryopType,
7   idToken,
8   biexpType,
9   keywords,
10  typeToken,
11  intToken,
12  epsilon
13 } nodeTypeEnum;

```

Besides the non-terminals given by the guidance, here some other AST types are introduced:

- a. Terminal types: keywords, typeToken, intToken, idToken, epsilon
- b. Auxiliaries: unaryopType, biexpType

”Auxiliaries” is special cases of *EXPS* non-terminal. *unaryopType* is for unary operators, while *biexpType* is for all binary expressions other than those concerned with assignment.

Since the sub AST varies from AST node types, a variable arguments functions is implemented to create AST node in the process.

```

1 TreeNode* createNode(int lineno, nodeTypeEnum type, char* data, int
   cnt, ...) {
2   TreeNode* ptr = (TreeNode*) malloc(sizeof(TreeNode));
3   ptr->lineNum = lineno;
4   ptr->type = type;
5   ptr->data = strdup(data);
6   ptr->size = ptr->capacity = cnt;
7   ptr->children = (TreeNode**) malloc(sizeof(TreeNode*) * cnt);
8   va_list ap;
9   va_start(ap, cnt);
10  int i;
11  for (i = 0; i < cnt; i++) {
12      ptr->children[i] = va_arg(ap, TreeNode*);
13  }
14  va_end(ap);
15  return ptr;
16 }

```

Then, the grammars of smallC is translated into *bison* file. Taken *EXTDEF* as instance:

```

1 EXTDEF      : TYPE EXT_VARS SEMI      { $$ = createNode(ylineno,
   extdefType, "extdef1", 2, createNode(ylineno, typeToken, $1, 0)
   , $2); }
2      | STSPEC EXT_VARS SEMI      { $$ = createNode(ylineno,
   extdefType, "extdef2", 2, $1, $2); }
3      | TYPE_FUNC SIMIBLOCK      { $$ = createNode(ylineno,
   extdefType, "extdef func", 3, createNode(ylineno, typeToken, $1
   , 0), $2, $3); }

```

```
4 ;
```

As the guidance indicated, grammar given may induce one or two reduce-reduce/shift-reduce conflicts. Specifically, *"IF LP EXP RP STMT* should have lower precedence than *IF LP EXP RP STMT ELSE STMT*. A trick of *bison* is applied here:

```
1 %nonassoc LOWER_THAN_ELSE
2 %nonassoc ELSE
3 %%
4 STMT: // .....
5     | IF LP EXPS RP STMT %prec IFX    {$$ = createNode(yylineno,
6       stmtType, "if stmt", 2, $3,$5); }
7     | IF LP EXPS RP STMT ELSE STMT %prec ELSE {$$ = createNode(
8       yylineno,stmtType, "if stmt", 3, $3,$5,$7);}
9 ;
```

One more modification is applied to grammar. As *read* and *write* is regarded as keywords in the project, the production of *STMT* now is

```
1 STMT      : EXP SEMI          {$$ = ...; }
2           | STMTBLOCK        {$$ = ...; }
3           | RETURN EXPS SEMI  {$$ = ...; }
4           | IF LP EXPS RP STMT %prec IFX    {$$ = ...; }
5           | IF LP EXPS RP STMT ELSE STMT %prec ELSE {$$ = ...;}
6           | FOR LP EXP SEMI EXP SEMI EXP RP STMT  {$$ = ...; }
7           | CONT SEMI        {$$ = ...; }
8           | BREAK SEMI       {$$ = ...; }
9           | READ LP EXPS RP SEMI          {$$ = ...;}
10          | WRITE LP EXPS RP SEMI         {$$ = ...;}
```

4 Abstract Syntax Tree

This part is done in *ast.h* and *def.h*.

With *samlLC.y* finished, the program now can output the AST. The struct *TreeNode* has already implemented the tree, thus the generation of AST is quite simple with techniques learned from data structure classes.

```
1 void printAst(TreeNode* ptr, int depth) {
2     int i;
3     int n = (depth - 1) * 2;
4
5     if (depth > 0) {
6         for (i = 0; i < n; i++) {
7             putchar(buffer[i]);
8         }
9         putchar(' | ');
10        puts("");
11    }
12
13    for (i = 0; i < n; i++) {
14        putchar(buffer[i]);
15    }
```

```

16     if (depth > 0) {
17         putchar('|');
18         putchar(' ');
19     }
20     n = depth * 2;
21     buffer[n] = '|';
22     buffer[n + 1] = ' ';
23     if (ptr == NULL) {
24         printf("NULL\n");
25         return;
26     }
27     printf("%s\n", ptr->data);
28     // terminal customized configuration
29     if (ptr->size > 0) {
30         printf("%s\n", ptr->data);
31     } else {
32         printf("%s\n", ptr->data);
33
34         for (i = 0; i < n; i++) {
35             putchar(buffer[i]);
36         }
37         puts("");
38     }
39     for (i = 0; i < ptr->size; i++) {
40         if (i == ptr->size - 1) {
41             buffer[n] = ' ';
42         }
43         printAst(ptr->children[i], depth + 1);
44     }
45 }
46

```

The AST of test.c is stored at file *AST*. It turned out that the syntax analyzer works well.

5 semantic analysis

This part is done in *semantic.h*.

In this part, potential semantic errors after generating a parse tree is checked. Besides, the symbol table is generated. Some basic evaluation (i.e. get values of int) is also done in this part.

5.1 Symbol Table

In this part, a struct is designed to implement the symbol table.

```

1 struct ptr_cmp {
2     bool operator()(const char* s1, const char* s2) const {
3         return strcmp(s1, s2) < 0;
4     }
5 };
6 struct SymbolTable {
7     map<char*, char*, ptr_cmp> table;
8     map<char*, vector<char*>, ptr_cmp> struct_table;

```



```

9 | map <char*, vector<char *>, ptr_cmp > struct_id_table;
10 | map < char*, int, ptr_cmp > struct_name_width_table;
11 |     map <char*, int, ptr_cmp> width_table;
12 | map <char*, vector<int>, ptr_cmp> array_size_table;
13 | int parent_index;
14 | } symTable[MAX_LEVEL][MAX_DEPTH];

```

Here STL data structure *map* is used for string-key inquiry.

As the member variable name indicates, they are used to map between:

member	string	key
table	local/global variable names	variable's type
struct_table	struct names	all member variable names of the struct
struct_id_table	struct names	member names certain type of the struct
struct_name_width_table	struct names	the number of members of the struct
width_table	int/array names	space it consumes
array_size_table	array names	how many 4 bytes of the width each dimension

Table 1: Symbol Table Content

symTable[MAX_LEVEL][MAX_DEPTH]: Here the symbol table is implemented as 2D array. The row index indicates the level of certain table. The smallC code has code blocks(" "), which requires a new symbol table. when the program is running, the symbol table of each code blocks has its own scope. Thus it is necessary to keep the level on track.

Besides, member variable *parent_index* is used to record the upper level's index. *parent_index* will be 1 for global symbol table.

5.2 Semantic Checking

For each AST node type, a corresponding check function is implemented. Normally it will get the AST node type of its sub-AST and call corresponding check function to check its sub-tree one by one. i.e.

```

1 | void check_program(TreeNode* p) {
2 |     if (p == NULL) return;
3 |     int i;
4 |     for (i = 0; i < p->size; ++i) {
5 |         check(p->children[i]);
6 |     }
7 | }

```

Specifically, the following things need to be checked: (since the amount of code in this part is really large, only design logic instead of code is presented here)

5.2.1 Variables and functions should be declared before usage

For variables, every time an identifier node is encountered, the compiler will check it in the symbol table on current level and its upper levels.

For functions, the compiler only allows global functions, which is checked on level 0 symbol table. *Overloaded function* is allowed in the compiler, thus not only the function name but its arguments will be checked in this part.

5.2.2 Variables and functions should not be redeclared

Similarly to the previous one, whenever a new variable is declared, it will be checked in the current-level symbol table to find out whether it has already existed.

5.2.3 Reserved words can not be used as identifiers

Whenever an identifier node is encountered, its name(stored in the *data*) will be compared to keywords.

5.2.4 Program must contain a function `int main()` to be the entrance

A global variable *bool_main* is initialized to be false. Once *main* is detected in *FUNC* node, it will be true, otherwise the semantic checking won't pass.

5.2.5 The number and type of variable(s) passed should match the definition of the function

This part is already done as a by-product of part 5.2.1

5.2.6 Use `[]` operator to a non-array variable is not allowed

With *width_table* from symbol table, the compiler can get access to the space an array or variable consumes. From it the compiler can tell whether it is an array or a int, then `[]` operation to int will be rejected.

Note. The declaration like `a[1]` is not allowed by the compiler, since it cannot distinguish it with ordinary int.

5.2.7 The `.` operator can only be used to a struct variable

With *struct_table* from symbol table, the compiler can check whether a variable has been declared as struct. Then with *struct_id_table*, the identifier after `.` is checked whether it is a member of the struct.

5.2.8 *break* and *continue* can only be used in a *for*-loop

A global bool variable is declared to tell whether the current *STMT* node is inside a loop.

5.2.9 Right-value can not be assigned by any value or expression

This part will be implemented in the intermediate representation section.

5.2.10 The condition of if statement should be an expression with int type

The logic is similar to part 5.2.6, *width_table* is checked to tell the type of expression.

5.2.11 The condition of for should be an expression with int type or epsilon

5.2.12 Only expression with type int can be involved in arithmetic

These parts are similar to 5.2.10

5.2.13 Other details

- a. Size of array cant be negative or 0
- b. Members of a struct and argument of a function cant have same names

6 Intermediate Representation Code Translation

This part is done in *translate.h*.

Here comes the most tough, complex and exhausting part. After passing semantic analysis, the smallC code can be regarded as a correct one and the AST node will be translated into intermediate representation(IR).

6.1 Intermediate Representation

Here a MIPS-like three-address code is selected as IR. The different between IR and MIPS is IR has infinite registers so that every variable can be hold in the registers.

Like before, a struct is designed to implement IR instructions.

```
1 // for MIPS generation
2 struct Address{
3     RegType type; //type of the argument
4     string name; //For labels and function calls
5     int value; // register value
6     int real; // physical register in MIPS
7     int needload; // whether load op is necessary
8     int needclear; // whether this op should be cleared
9 };
10
11 struct Quadruple{
12     string op; // name of operation
13     int active; //whether it's active
14     int flag; //whether its arguments should be modified
15     Address arguments[3]; //the three arguments
16 };
```

RegType in struct Address is an enum type, which represent the type of corresponding address. It can be:

- a. ADDRESS_LABEL: The address is jump label number
- b. ADDRESS_CONSTANT: The address is constant number
- c. ADDRESS_TEMP: It represents a register
- d. ADDRESS_NAME: It represents a name

opname	Destination	Operand 1	Operand 2	Description
or	a	b	c	$a = b c$
xor	a	b	c	$a = b\hat{c}$
and	a	b	c	$a = b\&c$
sll	a	b	c	$a = b \ll c$
srl	a	b	c	$a = b \gg c$
add	a	b	c	$a = b + c$
sub	a	b	c	$a = b - c$
mul	a	b	c	$a = b * c$
div	a	b	c	$a = b/c$
rem	a	b	c	$a = b\%c$
neg	a			$a = -b$
lnot	a	b		$a = !b$
not	a	b		$a = \sim b$
beqz		a	label	if a==0 goto label
bnez		a	label	if a!=0 goto label
bgez		a	label	if a \geq 0 goto label
bgtz		a	label	if a \geq 0 goto label
blez		a	label	if a \leq 0 goto label
bltz		a	label	if a \leq 0 goto label
li	a	b(constant)		$a = b$
lw	a	b	c(constant)	$a = \text{mem}[b+c]$
sw	[memory]	a,b,c(c is constant)		$\text{mem}[b+c] = a$
move	a	b	c	$a = b$
label		l		l is int, set label
goto		l		l is int, unconditional jump
func		s		s is string, set function
call		s		s is string, call function

Table 2: Three Address Code Logic

Some vectors are also used to help the translation process:

```

1 vector <Quadruple> IR, GIR, MIR;
2 vector <int> RegisterState, RegisterOffset, LabelCont, LabelBreak,
   vs_reg;
3 vector <string> vs_id;
```

- a. RegisterState: The register can either hold value or address. RegisterState is used to tell its current state
- b. RegisterOffset: contains the info of offsets of the register.
- c. LabelBreak and LabelCont: indicates the label for break and continue operation.
- d. vs_reg and vs_id: associates identifiers with the registers.

6.2 Translate

6.2.1 EXPS translate

For *EXPS* AST ndoe, the compiler will classify it into following cases:

Arithmetic A new register will be allocated to store the result of the arithmetic process

Logical The following IR is used to assign the bool value to register:

```

1 li [register],1
2 beqz/bnez/bgez/bgtz/blez/bltz judgement label
3 li [register],0
4 label 1

```

Unary For ordinary ones, a new register will be allocated.

For "+" and "-", not only a new register is allocated but the one participating in the operation will be modified as well.

Assignment First, the operations like "*" will be decomposed into arithmetic one and pure assignment.

Then for pure assignment, the compiler will just do register modifications.

Array and Struct Offset will be applied to visit the member.

6.2.2 FUNC translate

The arguments of functions are stored on the stack.

Every time a function(except main()) call is encountered, the \$ra will be stored in the stack. The stack pointer will be restalled after function call.

6.2.3 STMT translate

if goto label in IR can be used here. That is if the judgement of if is false, the compiler will use jump to else statement(let's call it s2); otherwise it will first execute the statement right after judgement(let's call it s1). Then a *goto* is placed right after s1 to jump out of the block, which will not execute s2 in the case.

for *goto* label in IR can be used here, too. The basic logic is almost the same except there are one more jump label to jump back in order to form a loop.

read and write

The *read()* is to assign the value of register \$v0 into the register associated with variable, while *write()* assign value to register \$a0.

It is implemented with *scanf_one* and *printf_one*

```
1  __printf_one :  
2  li $v0, 1  
3  syscall  
4  jr $ra  
5  __scanf_one :  
6  li $v0, 5  
7  syscall  
8  jr $ra
```

7 Optimization

Subexpression elimination and dead code elimination is done in this part.

7.1 Subexpression elimination

This part is done in *optimize.h*

Subexpression elimination are like change:

```
1  move t2, t1  
2  move t3, t2
```

into

```
1  move t3, t1
```

and

```
1  li t1, 1  
2  move t2, t1
```

into

```
1  move t2, 1
```

Here t2 must not be reassigned afterwards. A flag is used to indicate whether the register is reassigned later.

7.2 Instruction Selection

This part is done in *interpret.h*

The following operation is done to reinforcement optimization:

- a. Check labels and delete those whose destination immediately following.
- b. For all arithmetic expressions, interpret into intermediate instructions.
- c. Check for successive lws and sws. Delete the sw where the the former instruction is lw and their target are the same.
- d. use shift operations when multiplying 2 or 4.

8 Machine Code Generation

This part is done in *interpret.h* and *optimize.h*

After optimization, the compiler will interpret the IR into MIPS code. The major difference between IR and MIPS is the number of registers. Thus the compiler shall implement register allocation.

Here the compiler will scan through \$11 to \$25 to find out the free register. If no free one, it will scan \$8 to \$10. If no register is available currently, *lw* and *sw* will be applied to store variables into memory.

8.1 Project Test

The compiler is already finished, the test is necessary.

The test.c is translated into

```
1 .globl main
2
3 .data
4 temp: .space 44
5
6 .text
7 __printf_one:
8 li $v0, 1
9 syscall
10 jr $ra
11 __scanf_one:
12 li $v0, 5
13 syscall
14 jr $ra
15
16 main:
17 la $a3, temp
18 add $11, $29, -40
19 move $13, $11
20 move $11, $13
21 jal __scanf_one
22 sw $2, 0($11)
23 move $11, $13
```

```

24 move $12, $13
25 lw $12, 0($12)
26 add $12, $12, 10
27 sw $12, 0($11)
28 move $11, $13
29 lw $11, 0($11)
30 move $4, $11
31 jal __printf_one
32 li $11, 0
33 add $12, $29, 4
34 sw $11, 0($12)
35 j __program_end
36 j __program_end
37 __program_end:
38 li $v0, 10
39 syscall

```

Data	Text	
Text		
[00400008]	24a60004	addiu \$6, \$5, 4
[0040000c]	00041080	sll \$2, \$4, 2
[00400010]	00c23021	addu \$6, \$6, \$2
[00400014]	0c10000f	jal 0x0040003c [main]
[00400018]	00000000	nop
[0040001c]	3402000a	ori \$2, \$0, 10
[00400020]	0000000c	syscall
[00400024]	34020001	ori \$2, \$0, 1
[00400028]	0000000c	syscall
[0040002c]	03e00008	jr \$31
[00400030]	34020005	ori \$2, \$0, 5
[00400034]	0000000c	syscall
[00400038]	03e00008	jr \$31
[0040003c]	3c071001	lui \$7, 4097 [temp]
[00400040]	23abffd8	addi \$11, \$29, -40

Figure 3: Test on SPIM

It runs on QtSPIM precisely. The compiler works well.

9 Acknowledgement

This project is large and tough. After finishing it, I learned a lot about compilers, C language and MIPS.

I would like to express my sincere thanks to instructor Prof. Jiang and teaching assistant. The class and project taught me a lot. Thank you for your contribution!