

COMP3109

Assignment 3

Writing a Compiler – Group Assignment (10 marks)

This third assignment is a **group assignment**, and is due **on Friday Week 11, 5pm**. Your assignment will not be assessed unless all two of the following criteria are met:

1. Hand in a signed group academic honesty form in the tutorial
2. For this assignment write at least 15 test cases and submit it with your assignment
3. Submit a tarball of your source code to elearning for all solved tasks with plenty of remarks, i.e., documentation should be in form of remarks.

Please form groups of two or three students (groups of three preferred). Your code must run on the `ucpu[01]` machines. When submitting, please provide a `README` file instructing the marker how to run your code.

Task 1

(10 marks)

We want to implement a mini compiler for a vector language that translates input programs to Intel's assembly code using SSE extensions. The compiler should be implemented in a language tool of your choice that is able to deal with attribute grammars (e.g. ANTLR, YACC, etc; note that we recommend ANTLR). A module in the vector language consists of several functions. Each function has several input and output vectors. We further assume that all vectors of a function have the same length. A function in the vector language consists of a vector assignment, sequences of vector assignments, if- and while-statements, and calls to other functions defined in the vector language.

The grammar of the language is given below, where M is the start symbol.

$M \rightarrow F M$	
$M \rightarrow \varepsilon$	$E \rightarrow E + E$
$F \rightarrow \text{func ident } P D S \text{ end}$	$E \rightarrow E - E$
$P \rightarrow (L)$	$E \rightarrow E * E$
$L \rightarrow \text{ident}$	$E \rightarrow E / E$
$L \rightarrow \text{ident} , L$	$E \rightarrow \text{min} (E , E)$
$D \rightarrow \text{var } L ;$	$E \rightarrow (E)$
$D \rightarrow \varepsilon$	$E \rightarrow \text{ident}$
$S \rightarrow \text{if } C \text{ then } S \text{ else } S \text{ endif}$	$E \rightarrow \text{num}$
$S \rightarrow \text{while } C \text{ do } S \text{ endwhile}$	$C \rightarrow E < \text{num}$
$S \rightarrow S ; S$	$\text{ident} \rightarrow [a-zA-Z_][a-zA-Z0-9_]*$
$S \rightarrow \text{ident} = E$	$\text{num} \rightarrow [0-9]+(\.[0-9]+)?$
$S \rightarrow \varepsilon$	

Write an attribute grammar that translates an input program in the vector language to an assembly file. The assembly output should take advantage of the SSE extension (high speed floating point operations which can operate on 4 single precision numbers at the same time) of the x86-64 architecture. The produced assembly code should define functions that can be used in C programs. The structure of this exercise is to put together assembly templates (i.e. blocks of text) which are given in the rest of this assignment. These text blocks need to be pasted and written to a file. Some of these text blocks need specific parameters that are to be computed by the attributes of the grammar.

For example the following module in the vector language defines a function **mymin**

```
1  func mymin(a, b, c)
2    c = min(a,b) + 20
3  end
```

that has three parameters **a**, **b** and **c**. The element-wise minimum is computed, the constant vector 20 is added, and the result is stored in parameter **c**. The function **mymin** can be used as sketched in the following C code fragment:

```
1  #include <stdlib.h>
2
3  /* alignment macro: aligns a memory block a to multiplies of a */
4  #define align(s,a) (((size_t) (s) + ((a) - 1)) & ~((size_t) (a) - 1))
5  /* Alignment for SSE unit */
6  #define SSE_ALIGN (16)
7  /* Number of elements */
8  #define NUM (100)
9
10 extern void mymin(long, float *, float *, float *);
11
12 int
13 main(void) {
14     float *a = malloc(sizeof(float) *NUM + SSE_ALIGN),
15           *b = malloc(sizeof(float) *NUM + SSE_ALIGN),
16           *c = malloc(sizeof(float) *NUM + SSE_ALIGN);
17     /* make sure that pointers are aligned to multiplies of 16 bytes */
18     a = (float *) align(a, SSE_ALIGN);
19     b = (float *) align(b, SSE_ALIGN);
20     c = (float *) align(c, SSE_ALIGN);
21     ...
22     /* write values to a and b */
23     ...
24     /* invoke the function written in the vector language */
25     mymin(NUM, a, b, c);
26     ...
27     /* read values from c */
28     ...
29     return 0;
30 }
```

In this code fragment three single precision floating point arrays **a**, **b** and **c** are declared and used as parameters of the vector language function **mymin**. Before using the arrays as vectors they need to

be aligned to multiples of 16 – this is an SSE requirement. Additionally, we assume that the length of a vector is a multiple of 4 because a single SSE instruction can operate on 4 single precision floating point numbers at the same time. The first parameter gives the length of the vectors followed by the vector parameters. Depending on the semantics of the function, the function either reads or writes from the arrays.

When you write assembly code you need to be aware of the processor architecture and how the program stack is structured (see textbook) so that you can access the length of the vectors and vector parameters. The 64bit extension of the Intel's 32bit architecture is called *x86-64* and it has some differences to the standard Intel 32bit architecture. The *x86-64* architecture has following features:

- We have sixteen 64bit-registers called `%rax`, `%rbx`, `%rcx`, `%rdx`, `%rsi`, `%rdi`, `%rbp`, `%rsp`, `%r8`, and `%r9` - `%r15`. The register `%rsp` is a special register that holds the top of the stack pointer and the register `%rbp` is normally used as a frame pointer register.
- Memory addresses are 64bit long.
- The first six integer and pointer parameters are passed via registers rather than the program stack. The seventh parameter and any following parameters after the seventh parameter are passed via the program stack. Refer to following table:

Argument	Register
1.	<code>%rdi</code>
2.	<code>%rsi</code>
3.	<code>%rdx</code>
4.	<code>%rcx</code>
5.	<code>%r8</code>
6.	<code>%r9</code>

- Register `%rax` stores the return value of a function.
- Register `%rbx` and `%r12` - `%r15` are callee-saved, i.e., if a function modifies these registers it needs to restore them before the function terminates.
- Register `%xmm0` - `%xmm15` are floating point registers that can store four floating point numbers simultaneously.

Recall that the program stack grows downwards. For the *x86-64* architecture, the stack layout is given below,

...	higher addresses	
return address	8	<i>saved by the call instructions</i>
frame pointer of calling functions	0	<i>saved as soon as the function is invoked</i>
first free element on stack	-8	<i>value of frame pointer</i>
...	lower addresses	

The stack pointer in *x86-64* is denoted by `%rsp` and the frame pointer by `%rbp`. A function definition in the vector language uses following template to produce assembly code:

```

1  .text
2  .global <name>
3  .type <name>, @function
4  .p2align 4,,15
5
6  <name>:
7      # save current frame pointer on stack
8      pushq   %rbp
9      # set frame pointer
10     movq    %rsp, %rbp
11     # save callee-save registers that are used on stack
12     pushq   %rbx
13
14     <allocate memory for local vector variables>
15
16     <Insert the body of function here>
17
18     # epilog of a function
19     popq    %rbx      # restore reg %rbx
20     leave   # restore frame pointer
21     ret      # return

```

where `.text` denotes that the machine code should be stored in the text segment (see textbook) of the memory. The pseudo-mnemonic `.global <name>` defines a global name `<name>` for the linker, and the pseudo-mnemonic `.type` declares that the symbol `<name>` is a function. With `<name>:` you define an address for the function entry point. The next instruction pushes the frame pointer of the calling function onto the stack. The second instruction moves the current value of the stack pointer to the frame pointer. The third instruction saves register `%rbx` onto to stack since it is callee save, i.e., the callee expects that `%rbx` does not change its value. Note we will use `%rbx` for the code templates and hence it needs to be saved at the entry and restored at the exit.

After the first three instructions you add the code of reserving space on the stack for local variables and the assignments of the vector function. A function definition ends with loading the previous value of `%rbx` from the stack, and executing `leave`, which restores the frame pointer and the stack. Instruction `ret` loads the return address into the instruction pointer of the CPU and returns to the next instruction of the caller.

Local variables in the vector language are stored on the stack from $-16(\%rbp)$ onward. The memory block for all local variables are created by subtracting the product of the length of a vector and the number of local stack variables from the stack pointer and adding a padding for pointer alignment, i.e.,

```

1  # allocate <NUM> local variables
2  movq   %rdi, %rax
3  imulq  $4, %rax, %rax
4  addq   $16, %rax
5  imulq  $<NUM>, %rax, %rax
6  subq   %rax, %rsp
7  andq   $-16, %rsp

```

where `<NUM>` is the number of local variables. The number of local variables can be computed by an

attribute in your attribute grammar. This text block needs to be generated as a first text block in the function body (if local vector variables are in the function).

The following three templates are concerned about computing addresses of vector parameters, local variables, and constants. These three text blocks will be used for address computations in templates that either compute new vectors or assign a vector/constant to another vector. The address of a vector parameter is stored in an argument register. We assume that there are at most 5 vector parameters for a function since we have only six argument registers and the first one is used for storing the length of the vectors.

```
1 # place address of <N>th parameter into <destreg>
2 movq    <argreg-N+1>, <destreg>
```

where $\langle \text{argreg-N+1} \rangle$ is the $(n + 1)$ 'th argument register as listed above. The destination register $\langle \text{destreg} \rangle$ is one of the following registers `%rax`, `%r10` or `%r11` that are used in the code templates for vector assignment and vector operations. An address of a local vector variable is computed by

```
1 # place address of <N>th local variable into <destreg>
2 movq    %rdi, <destreg>
3 imulq   $4, <destreg>, <destreg>
4 addq    $16, <destreg>
5 imulq   $<N>, <destreg>, <destreg>
6 subq    %rbp, <destreg>
7 negq    <destreg>
8 andq    $-16, <destreg>
```

where $\langle N \rangle$ is the n th local variable (starting from 1) of the vector function. The address has to be computed at runtime because the length of the vectors is not known ahead of time. The last template loads the address of a constant $\langle X \rangle$ into a destination register:

```
1 # place address of <X> into <destreg>
2 movq    $.const<X>, <destreg>
```

where the label `.const<X>` is a new label. For generating constants we need a postamble at the end of the assembly code that describes the constant:

```
1 .data
2 .align 16
3 .const<X>:
4     .float <X>
5     .float <X>
6     .float <X>
7     .float <X>
```

where $\langle X \rangle$ is the floating point number x . We need four repetitions because an SSE operation does four operations in a single step.

In the following we introduce assembly code templates for basic forms. In the following, we will introduce the code templates for the basic forms. There are two types of basic forms as shown below,

ident = factor

or

ident = factor **op** factor

where “factor” is either a variable or a constant. Conceptually, you break up the right-hand side of an expression in basic forms. In this conceptual transformation (which is also known as linearization of the code), you will introduce numerous local variables, e.g., the input statement

```
x = a + b + c + d
```

will be conceptually transformed to

```
x1 = a + b;
x2 = x1 + c;
x3 = x2 + d;
```

Note that I do not suggest to perform this transformation. I suggest to create local variables on the fly by using attributes and an expression node in the syntax tree representing a basic form. Furthermore, you want to reuse local variables as much as you can too avoid stack size explosion.

The basic form “**ident** = factor” is translated to assembly with following template:

```

1      <load source address into %rax>
2      <load destination address into %r10>
3
4      movq    %rdi, %rbx      # load vector length into counter %rbx
5      shrq    $2, %rbx       # divide counter reg by 4
6                               # (per loop iteration 4 floats)
7      jz      .loop_end<X>   # check whether number is equal to zero
8
9      .loop_begin<X>:        # loop header
10
11     movaps   (%rax), %xmm0   # load source into %xmm0
12     movaps   %xmm0, (%r10)   # store %xmm0
13
14     # IMPORTANT: remove the following line only if %rax is
15     # pointing to a constant
16     addq     $16, %rax       # increment source pointer by (4 x float)
17
18     addq     $16, %r10       # increment destination pointer by (4 x float)
19
20     decq     %rbx            # decrement counter
21     jnz      .loop_begin<X> # jump to loop header if counter is not zero
22
23     .loop_end<X>:
```

where <X> is a unique number for this basic form. This text block loads the source and destination address into %rax and %r10 by using the templates as previously shown. The assembly code inside loads the vector length and divides it by 4 (i.e. shift of number by two digits). Inside the loop the **movaps** instructions load the instructions into the first SSE register. After performing the move instruction, the first SSE register will contain 4 consecutive floating point numbers. With the second **movaps** instruction the 4 consecutive floating point numbers are written to the destination.

The basic form **ident** = factor **op** factor is translated with following template:

```

1    <load source1 address into %rax>
2    <load source2 address into %r10>
3    <load destination address into %r11>
4
5    movq    %rdi, %rbx    # load vector length into counter %rbx
6    shrq    $2, %rbx     # divide counter reg by 4
7                                # (per loop iteration 4 floats)
8    jz      .loop_end<X> # check whether number is equal to zero
9
10   .loop_begin<X>:      # loop header
11
12   movaps   (%rax), %xmm0 # load first operand into %xmm0
13   movaps   (%r10), %xmm1 # load second operand into %xmm1
14
15   # perform operation
16   <operation> %xmm1, %xmm0
17
18   movaps   %xmm0, (%r11) # store result
19
20   # increment pointers
21
22   # IMPORTANT: remove following line if %rax is pointing to a constant
23   addq     $16, %rax
24
25   # IMPORTANT: remove following line if %r10 is pointing to a constant
26   addq     $16, %r10
27
28   addq     $16, %r11
29
30   decq     %rbx          # decrement counter
31   jnz      .loop_begin<X> # jump to loop header if counter is not zero
32   .loop_end<X>:

```

Use the operation **addps** for addition, **subps** for subtracting, **divps** for division, **mulps** for multiplication, and **minps** for minimum.

Make sure that the code generation works for sequence assignments first before you attempt if-statements and while-loops. If-statements and while-loops have a condition of the form $E < \text{num}$ that computes the sum of the vector elements of E and checks whether the sum is less than num . Depending on the outcome either the true-branch or false-branch is executed or the while loop terminated. For the generation of assembly code for if-statements and while-loops you need templates that generate the necessary jump instructions.

For if-statements the template is the following

```

1    <template for condition(.true-branch<num> , .false-branch<num>)
2
3    .truebranch<num>:
4
5    <emit code for true-branch here>
6
7    jmp .endif<num>

```

```

8
9 .falsebranch<num>:
10
11 <emit code for false-branch here>
12
13 .endif<num>:

```

where for each if-statement you need unique labels. The template for a while loop is given below

```

1  jmp .loopcond<num>
2
3 .loopbegin<num>:
4  <emit code for loop-body here>
5
6 .loopcond<num>:
7  <template for condition (.loopbegin<num>, .loopexit<num>)>
8
9 .loopexit<num>:

```

where for each while-statement you need unique labels as well.

The templates for the sums are given below:

```

1  <load source address into %rax>
2
3  xorps %xmm1, %xmm1 # set register %xmm1 to zero
4  movq  %rdi, %rbx   # load vector length into counter %rbx
5  shrq  $2, %rbx     # divide counter reg by 4
6                      # (per loop iteration 4 floats)
7  jz    .loop_endX>  # check whether number is equal to zero
8
9  .loop_beginX>:      # loop header
10
11
12  addps (%rax), %xmm1 # add source to %xmm0
13
14  # IMPORTANT: remove the following line only if %rax is
15  # pointing to a constant
16  addq  $16, %rax     # increment source pointer by (4 x float)
17
18  decq  %rbx          # decrement counter
19  jnz   .loop_beginX> # jump to loop header if counter is not zero
20
21 .loop_endX>:
22  xorps %xmm0, %xmm0 # set register %xmm0 to zero
23  addss %xmm1, %xmm0  # add all four numbers in %xmm1 to %xmm0
24  shufps $147, %xmm1, %xmm1 # note that
25  addss %xmm1, %xmm0  # a shuffle operation rotates the single precision
26  shufps $147, %xmm1, %xmm1 # to the left. In each step
27  addss %xmm1, %xmm0  # a number is read and added to the single
28  shufps $147, %xmm1, %xmm1 # precision number in %xmm0
29  addss %xmm1, %xmm0

```



```
30
31                                     # compare sum with number and jump
32                                     # to the true/false target
33     ucomiss .L<number>, %xmm0
34     ja     <true>
35
36     jmp     <false>
37
38     .align 4
39     .L<number>:
40     .float <number>
```

After you have written the stub in C and produced code you can create an executable with

```
gcc -Wall -W -g -o myexe stub.c compiled.s
```

and you can trace the executable with `gdb`. The emission of the assembly code might be done as a traversal in your parser generator tool to avoid awkward data structures for storing the text file.

For speed fanatics: you can achieve substantial more speed by unrolling the loop and pre-computing the addresses of local variables.

Two scripts which might come in use. This first script called `compile.sh`, is used to compile our sample solution to the final executable. Our ANTLR grammar file is called `VPL.g` (Vector Processing Language).

```
1  #!/bin/bash
2  set -e
3
4  BUILD_DIR=build
5
6  if [ $# != 1 ]; then
7      echo "Usage: ${0} filename.vpl" >&2
8      exit 1
9  fi
10
11 # builds the ANTLR-generated parser using the grammar file
12 java -cp antlr-3.1.2.jar org.antlr.Tool -o ${BUILD_DIR} VPL.g
13 touch ${BUILD_DIR}/__init__.py
14
15 # uses the ANTLR-generated parser to convert the VPL program to ASM
16 ./vpl2asm.py < ${1} > ${1}.s
17
18 # compiles the ASM and C file together
19 gcc -Wall -W main.c ${1}.s -o my_program
```

The second script, `vpl2asm.py`, invokes our Python parser generated by ANTLR:

```
1  #!/usr/bin/env python
2  import sys
3  import antlr3
4  from build.VPLLexer import VPLLexer
5  from build.VPLParser import VPLParser
6
7  char_stream = antlr3.ANTLRInputStream(sys.stdin)
8  lexer = VPLLexer(char_stream)
9  tokens = antlr3.CommonTokenStream(lexer)
10 parser = VPLParser(tokens)
11 root = parser.prog()
```

These scripts are useful if you use Python as your target language within ANTLR. However, they can be easily adapted for your own chosen target language.