

CSC488
ASSIGNMENT 5
ASSEMBLER

DANIEL BLOEMENDAL

CONTENTS

1. Overview	1
2. Grammar	1
3. Labels	2
4. Instructions	2
4.1. PUSHSTR(\langle string_operand \rangle string)	2
4.2. SETUPCALL(\langle value_operand \rangle return_address)	2
4.3. JMP(\langle value_operand \rangle address)	2
4.4. BFALSE(\langle value_operand \rangle address)	2
4.5. NOT	2
4.6. SAVECTX(\langle value_operand \rangle lexical_level)	3
4.7. RESTORECTX(\langle value_operand \rangle lexical_level, \langle value_operand \rangle argument_count)	3
4.8. RESERVE(\langle value_operand \rangle words)	3
5. Example	3
6. Error handling	4
7. Changes	4
8. Credits	5

1. OVERVIEW

The first goal, before beginning work on code generation, was to implement an assembler for our IR instruction set. This included instructions with label operands, instead of static addresses, and each macro, or IR instruction, outlined in the code generation templates. The entire IR instruction set will be covered again in detail in this document. There are also a few notable bugs that were fixed since the code generation templates in A4.

The assembler was designed to be decoupled from the code generator. It is standalone and can even be executed from command line on IR assembly programs. A run script is provided along with an “ant” command to build the assembler. To build the assembler, run “ant assembler”. To execute the assembler from command line, run “./RUNASSEMBLER.sh <source.ir>”

2. GRAMMAR

The following is the grammar for IR assembly programs in EBNF.

Please refer to http://en.wikipedia.org/wiki/Extended_Backus-Naur_Form for more information.

$\langle \text{program} \rangle$::= { $\langle \text{section} \rangle$ }
$\langle \text{section} \rangle$::= 'SECTION', ',', $\langle \text{ident} \rangle$, '\n', { [$\langle \text{label} \rangle$], [$\langle \text{instruction} \rangle$], '\n' }
$\langle \text{label} \rangle$::= $\langle \text{ident} \rangle$, ':'
$\langle \text{instruction} \rangle$::= ($\langle \text{machine_operation} \rangle$ $\langle \text{ir_operation} \rangle$), { $\langle \text{operand} \rangle$ }
$\langle \text{machine_operation} \rangle$::= 'HALT' 'ADDR' 'LOAD' 'STORE' 'PUSH' 'PUSHMT' 'SETD' 'POP' 'POPNT' 'DUP' 'DUPN' 'BR' 'BF' 'NEG' 'ADD' 'SUB' 'MUL' 'DIV' 'EQ' 'LT' 'OR' 'SWAP' 'READC' 'PRINTC' 'READI' 'PRINTI' 'TRON' 'TROFF' 'ILIMIT'
$\langle \text{ir_operation} \rangle$::= 'PUSHSTR' 'SETUPCALL' 'JMP' 'BFALSE' 'NOT' 'SAVECTX' 'RESTORECTX' 'RESERVE'
$\langle \text{operand} \rangle$::= $\langle \text{value_operand} \rangle$ $\langle \text{boolean_operand} \rangle$ $\langle \text{string_operand} \rangle$
$\langle \text{value_operand} \rangle$::= $\langle \text{integer_operand} \rangle$ $\langle \text{label_operand} \rangle$
$\langle \text{integer_operand} \rangle$::= ['-'], $\langle \text{digit} \rangle$, { $\langle \text{digit} \rangle$ }
$\langle \text{boolean_operand} \rangle$::= '\$true' '\$false'
$\langle \text{label_operand} \rangle$::= $\langle \text{ident} \rangle$
$\langle \text{string_operand} \rangle$::= '"', { ? all characters ? }, '"'
$\langle \text{ident} \rangle$::= $\langle \text{letter} \rangle$, { $\langle \text{letter} \rangle$ $\langle \text{digit} \rangle$ }
$\langle \text{letter} \rangle$::= 'a'..'z' 'A'..'Z' '_'
$\langle \text{digit} \rangle$::= '0'..'9'

3. LABELS

The biggest difference between the IR assembly language and the assembly language built into **Machine** is the support for labels. Labels can be thought of as instructions of with zero size and no operation. The addresses of labels are stored as instructions are assembled. Afterwards, any label operand pointing to a valid label is set to resolve to the computed address of the label it points to.

4. INSTRUCTIONS

In addition to support for labels, the assembler also provides an extended instruction set. The following instructions outlined below are not available in **Machine**.

4.1. **PUSHSTR**(\langle string_operand \rangle string). This instruction adds **string** to the text constant pool, if it has not already been added, and obtains the address of the string in the pool. The address is then pushed on to the stack.

```
PUSH 'address of string in constant pool'
```

4.2. **SETUPCALL**(\langle value_operand \rangle return_address). **SETUPCALL** sets up a call, according to calling convention designed for the compiler, by pushing a place holder for the return value, and by pushing **return_address**, the return address for the call. The placeholder for the return value is set to be **Machine.UNDEFINED**, to support returning from a function or procedure without an explicit “**return**” or “**result**” statement. For a more detailed discussion on the calling convention, see the code generation documentation.

```
PUSH UNDEFINED
PUSH return_address
```

4.3. **JMP**(\langle value_operand \rangle address). **JMP** pushes **address** on to the stack and then performs a branch.

```
PUSH address
BR
```

4.4. **BFALSE**(\langle value_operand \rangle address). Like **JMP**, **BFALSE** is also a short hand for pushing an address and then branching. This time, the branch is the conditional **BF**.

```
PUSH address
BF
```

4.5. **NOT**. Unfortunately the machine has no built in method of negating boolean values. **NOT** provides this much needed functionality. It does so by comparing the top of the stack with **\$false**. To see this, suppose $top = \$true$, then $(\$true = \$false) = \$false$ on the other hand if $top = \$false$ then $(\$false = \$false) = \$true$

```
PUSH $false
EQ
```

4.6. **SAVECTX**($\langle \text{value_operand} \rangle$ lexical_level). **SAVECTX** is used at the beginning of a major scope to preserve any previous display set by a major scope at the same lexical level, a sibling scope. After preserving the display, it sets the display to the current address stored in the stack pointer. This is part of the calling convention and is explained in more detail in the code generation documentation.

```
ADDR lexical_level 0
PUSHMT
SETD lexical_level
```

4.7. **RESTORECTX**($\langle \text{value_operand} \rangle$ lexical_level , $\langle \text{value_operand} \rangle$ argument_count). **RESTORECTX** is used at the end of major scope restore the display saved by **SAVECTX**. It also pops argument_count arguments off of the stack. Again, for more details on the calling convention, refer to the code generation documentation.

```
SETD lexical_level
PUSH argument_count
POPNT
```

4.8. **RESERVE**($\langle \text{value_operand} \rangle$ words). **RESERVE** is used to reserve words of memory on the stack. It is current used in the prolog of a major scope to reserve memory for locals.

```
PUSH 0
PUSH words
DUPN
```

5. EXAMPLE

The example below is the currently used runtime library for the code generator. It contains only one function right now, **print**. It serves as a good demonstration of the IR assembly, using a number of the different IR instructions.

```
1 ; -----
2 ; Start of runtime library
3 ; -----
4     SECTION .library
5
6 print:
7     SAVECTX 0
8     PUSH 0
9 __print_start:
10    DUP
11    ADDR 0 -2
12    LOAD
13    LOAD
14    LT
15    PUSH __print_end
16    BF
17    DUP
18    ADDR 0 -2
```

```
19      LOAD
20      PUSH 1
21      ADD
22      ADD
23      LOAD
24      PRINTC
25      PUSH 1
26      ADD
27      PUSH __print_start
28      BR
29  __print_end: POP
30      RESTORECTX 0 1
31      BR
32  ; -----
33  ; End of runtime library
34  ; -----
```

6. ERROR HANDLING

The machine has a number of nontrivial limitations that are easy to run into. The limited memory and 16 bit operands ensure that it is easy to hit the limits. The assembler is careful to perform bounds checking on integer operands, enforcing the `Machine.MIN_INTEGER` and `Machine.MAX_INTEGER` bounds. In addition to integer bounds checking, the assembler also will gracefully handle programs that exceed the available machine memory `Machine.maxMemory`. An exception will be thrown when these errors are encountered.

7. CHANGES

Since assignment 4 bugs were identified and fixed. First, `RESTORECTX` incorrectly restored the display. Instead of loading the saved display via `ADDR`, and then popping it off the stack afterwards, it made more sense to simply consume it and restore the display simultaneously with `SETD`. Other issues included the correct offsets in a stack frame. The stack frame and calling conventions are discussed in detail in the code generation documentation.

8. CREDITS

src/compiler488/codegen/assembler/Operand.java	Daniel Bloemendal
src/compiler488/codegen/assembler/Section.java	Daniel Bloemendal
src/compiler488/codegen/assembler/StringOperand.java	Daniel Bloemendal
src/compiler488/codegen/assembler/Instruction.java	Daniel Bloemendal
src/compiler488/codegen/assembler/Processor.java	Daniel Bloemendal
src/compiler488/codegen/assembler/Assembler.java	Daniel Bloemendal
src/compiler488/codegen/assembler/InvalidInstructionException.java	Daniel Bloemendal
src/compiler488/codegen/assembler/IntegerOperand.java	Daniel Bloemendal
src/compiler488/codegen/assembler/ProgramSizeException.java	Daniel Bloemendal
src/compiler488/codegen/assembler/LabelNotResolvedException.java	Daniel Bloemendal
src/compiler488/codegen/assembler/LabelOperand.java	Daniel Bloemendal
src/compiler488/codegen/assembler/ir/AssemblerMachineEmitter.java	Mike Qin
src/compiler488/codegen/assembler/ir/AssemblerIREmitter.java	Mike Qin
src/compiler488/codegen/assembler/ir/Emitter.java	Mike Qin