

**CMSC 22610
Winter 2005**

**Implementation
of
Computer Languages**

**Project 4
March 1**

**MLR interpretation
Due: March 17**

1 Introduction

The final part of the project is to implement a code generator for MLR. The target for this code generator is a stack-based interpreter, which is described below.

2 The MLR virtual machine

In this section, we describe the MLR virtual machine (VM). The virtual machine is a stand-alone program that takes an executable file and runs it. An VM executable consists of a code sequence, a literal table that contains string literals, and a C function table that contains runtime system functions used to implement services such as I/O.

2.1 Values

The VM supports three types of values: 31-bit tagged integers, 32-bit pointers to heap-allocated records of values, and 32-bit pointers to strings. A integer value n is represented by $2n + 1$ in the VM (this tagging is required for the garbage collector). The VM takes care of tagging/untagging, so the only impact of this representation on your code generator is that integer literals must be in the range -2^{30} to $2^{30} - 1$. We use word address for values (but byte addressing for instructions).

2.2 Registers

The MLR VM has four special registers: the stack pointer (SP), which points to the current top of the stack; the frame pointer (FP), which points to the base of the current stack frame and is used to access local variables; the environment pointer (EP), which points to the current closure object and is used to access global variables; and the program counter (PC), which points to the next instruction to execute.

2.3 Instructions

We define the semantics of the instructions using the following notation

$$\dots \alpha \text{ instr} \implies \dots \beta$$

which means that the instruction **instr** takes a stack configuration with α on the top and maps it to a stack with β on the top. The instructions are organized by kind in the following description.

Arithmetic instructions

- $\dots i_1 i_2 \text{ add} \implies \dots (i_1 + i_2)$
pops the top two integers, adds them and pushes the result.
- $\dots i_1 i_2 \text{ sub} \implies \dots (i_1 - i_2)$
pops the top two integers, subtracts them and pushes the result.
- $\dots i_1 i_2 \text{ mul} \implies \dots (i_1 \times i_2)$
pops the top two integers, multiplies them and pushes the result.
- $\dots i_1 i_2 \text{ div} \implies \dots (i_1 / i_2)$
pops the top two integers, divides them, and pushes the result. The result is undefined if i_2 is zero.
- $\dots i_1 i_2 \text{ mod} \implies \dots i_1 \bmod i_2$
pops the top two integers, divides them, and pushes the remainder. The result is undefined if i_2 is zero.
- $\dots i \text{ neg} \implies \dots -i$
pops the integer on the top of the stack and pushes its negation.
- $\dots v_1 v_2 \text{ equ} \implies \dots b$
pops and compares the two values on top of the stack. If they are equal, then it pushes 1, otherwise it pushes 0. Note that if the values are pointers, then the comparison pushes true if the pointers are equal
- $\dots v_1 v_2 \text{ less} \implies \dots b$
pops and compares the two integers on top of the stack. If $v_1 < v_2$, then it pushes 1, otherwise it pushes 0.
- $\dots v_1 v_2 \text{ lesseq} \implies \dots b$
pops and compares the two integers on top of the stack. If $v_1 \leq v_2$, then it pushes 1, otherwise it pushes 0.
- $\dots v \text{ not} \implies \dots b$
pops v and pushes 1, if $v = 0$, and otherwise pushes 0.

Heap instructions

- $\dots v_0 \dots v_{n-1} \text{ alloc}(n) \implies \dots \langle v_0, \dots, v_{n-1} \rangle$
allocates an n element record, which is initialized from the top n stack values.
- $\dots \langle v_0, \dots, v_{n-1} \rangle \text{ select}(i) \implies \dots v_i$
pops a record off the stack and pushes the record's i th component.
- $\dots \langle v_1, \dots, v_{n-1} \rangle i \text{ index} \implies \dots v_i$
pops a record off the stack and pushes the record's i th component.
- $\dots \langle v_0, \dots, v_{n-1} \rangle i v \text{ update} \implies \dots$
replaces the i th field of the record with the value v .

Stack instructions

- ... **int**(n) \implies ... n
pushes the integer n onto the stack.
- ... **literal**(i) \implies ... s_i
pushes a reference to the i th string literal (s_i) onto the stack.
- ... **label**(l) \implies ... $addr$
pushes the code address named by the label. Note that in the encoding of this instruction, the code address is specified as an offset from the **label** instruction.
- ... $v_1 v_2$ **swap** \implies ... $v_2 v_1$
swaps the top two stack elements.
- ... $v_0 v_1 \dots v_{n-1} v_n$ **swap**(n) \implies ... $v_n v_1 \dots v_{n-1} v_0$
swaps the top stack element with the n 'th from the top. All other stack elements are unchanged.
- ... $v_n \dots v_0$ **push**(n) \implies ... $v_n \dots v_0 v_n$
pushes the n th element from the top of the stack.
- ... v **pop** \implies ...
pops and discards the top stack element.
- ... $v_1 \dots v_n$ **pop**(n) \implies ...
pops and discards the top n stack elements.
- ... **loadlocal**(n) \implies ... v
fetches the value (v) in the word addressed by $FP + n$ and pushes it on the stack.
- ... v **storelocal**(n) \implies ...
pops v off the stack and stores it in the word addressed by $FP + n$.
- ... **loadglobal**(n) \implies ... v
fetches the value (v) in the word addressed by $EP + n$ and pushes it on the stack.
- ... **pushep** \implies ... ep
push the current contents of the EP on the stack.
- ... ep **pop ep** \implies ...
pop a value from the stack and store it in the EP.

Control-flow instructions

- ... **jmp**(n) \implies ...
transfer control to instruction $PC + n$.
- ... b **jmpif**(n) \implies ...
pops b off the stack and if $b \neq 0$ it transfers control to instruction $PC + n$.
- ... **call** \implies ... pc
pushes the current PC value (which will be the address of the next instruction) and transfers control to $addr$, where the current value of the EP is $\langle addr, \dots \rangle$.

- ... **entry**(n) \implies ... $fp\ w_1 \dots w_n$
 pushes the current value of the FP register and sets FP to SP. Then it allocates n uninitialized words on the stack.
- ... $pc\ fp \dots v$ **ret** \implies ...
 resets the stack pointer to the frame-pointer; pops the saved FP into the FP register, pops the return PC, and then jumps to the return address.
- ... $pc\ fp \dots v$ **tailcall** \implies ... pc
 pops the current frame off the stack (like **ret**) and then transfers control to $addr$, where the current value of the EP is $\langle addr, \dots \rangle$. Unlike the **call** instruction, this instruction does not push the return PC.
- ... $args$ **ccall**(n) \implies ... v
 Calls the n th C function. The C function will pop its arguments from the stack and push its result.

Miscellaneous instructions

- ... **nop** \implies ...
 no operation.
- ... **halt** \implies ...
 halts the program.

3 Calling conventions

An MLR function call is implemented using a four-part protocol:

1. The caller evaluates the function and argument expressions from left to right and pushes the results onto the stack. Then a **swap** instruction is used to get the function closure on top of the stack. The closure is loaded into the EP register using the **popesp** instruction. Then the function is called (using the **call**) instruction, which has the effect of pushing the address of the following instruction on the stack.
2. The first instruction in the function is an **entry** instruction, which pushes the frame-pointer, sets the new frame pointer to point to the top of the stack, and then allocates space for locals.
3. When the callee is finished and the return result is on the top of the stack, it stores the result at the location of the argument (immediately below the return PC) and then executes a **ret** instruction, which deallocates the local variable space, restores the frame pointer, and transfers control to the return address.
4. When control is returned to the address following the **call**, the caller must save the return result, which will be on top of the stack.

There are two important variations on this protocol. The first is when a function calls itself. In that case, the EP already holds the closure and does not have to be set. The second case is when the function call is a *tail call*, i.e., the last action a function takes before returning. Tail calls are used to implement looping in functional languages. The VM has a special **tailcall** operator that discards the caller's stack frame and does not push the return PC.

3.1 An example

To illustrate the VM, consider the following MLR implementation of the factorial function:

```
fun fact (n : int) : int in
  if (n <= 0) then 1 else n * fact(n-1)
```

This function has one argument (*n*) and no local variables. The argument *n* will be located at +4 from the frame pointer, while *i* is at -2 and *p* is at -4 from the frame pointer. The VM code for this function is given in Figure 1 (we assume that the code is located at address 100). Consider a

Address	Instruction	Comment
100:	entry (0)	
102:	loadlocal (2)	push n
104:	int (0)	
106:	lesseq	is (n < 0)?
107:	jmpif (11)	jump to 120
109:	loadlocal (2)	push n
111:	loadlocal (2)	push n
113:	int (1)	
115:	sub	
116:	call	self-recursive call of fact
117:	mul	
118:	jmp (2)	jump to 122
120:	int (1)	
122:	storelocal (2)	save result in argument position
124:	ret	

Figure 1: AVM code for factorial function

non-tail call of the factorial function

```
fact x
```

This will produce the following code sequence:

```
loadlocal(fact)
loadlocal(x)
swap
popop
call
```

Here we have assumed that *fact* and *x* are local variables and we have used their names to refer to their offsets.

3.2 Submission

Your CVS repository will be seeded with CVS modules for each of the projects. For this project, the CVS module is named `project-4` and contains the sample scanner, parser, and typechecker implementation. We will also provide a code generation API. We will collect the projects at 12 midnight on Thursday March 17th from the repositories, so make sure that you have committed your final version before then.

4 Document history

Mar. 1 Original version.

Mar. 2 Added missing **update** operation.

Mar. 3 Changed the semantics of **ret** and **tailcall**, and added the **push** instruction.

Mar. 11 Changed documentation of **ccall** instruction.