

Appendix E

Definition of the CPRL Virtual Machine

E.1 Specification.

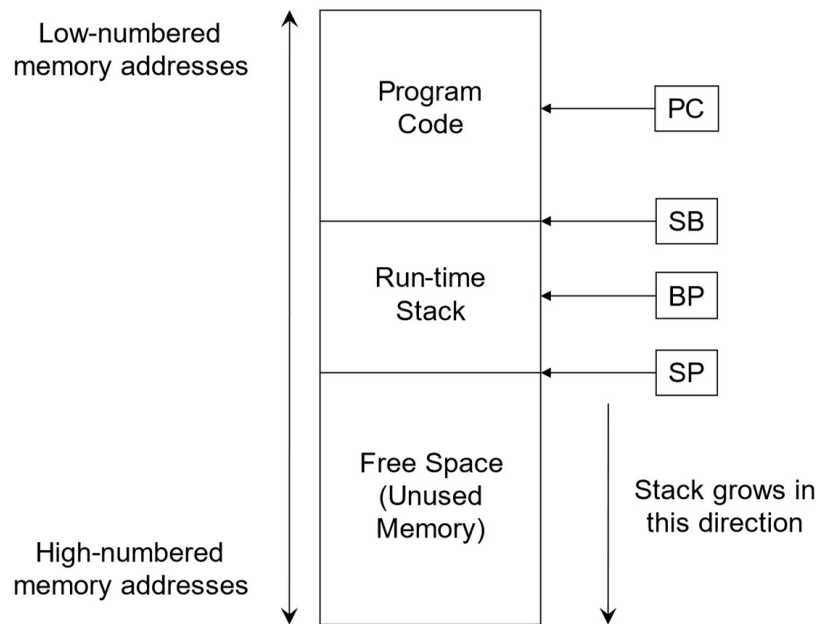
CVM (CPRL Virtual Machine) is a hypothetical computer designed to simplify the code generation phase of a compiler for CPRL (the Compiler PProject Language). CVM has a stack architecture; i.e., most instructions either expect operands on the stack, place results on the stack, or both. Memory is organized into 8-bit bytes, and each byte is directly addressable. A word is a logical grouping of 4 consecutive bytes in memory. The address of a word is the address of its first (low) byte. Boolean values are represented in a single byte, character values use 2 bytes (Unicode Basic Multilingual Plane or Plane 0, code points from U+0000 to U+FFFF), and integer values use a word (four bytes).

CVM has four 32-bit internal registers that are usually manipulated indirectly as a result of program execution. There are no general-purpose registers for computation. The names and functions of the internal registers are as follows.

- PC (program counter); a.k.a. instruction pointer: holds the address of the next instruction to be executed.
- SP (stack pointer): holds the address of the top of the stack. The stack grows from low-numbered memory addresses to high-numbered memory addresses. When the stack is empty, SP has a value of the address immediately before the first free byte in memory.
- SB (stack base): holds the address of the bottom of the stack. When a program is loaded, SB is initialized to the address of the first free byte in memory.
- BP (base pointer): holds the base address of the current activation record (a.k.a., frame); i.e., the base address for the subprogram currently being executed.

Each CVM instruction operation code (opcode) occupies one byte of memory. Some instructions take an immediate operand, which is always located immediately following the instruction in memory. Depending on the opcode, an immediate operand can be a single byte, two bytes (e.g., for a char), four bytes (e.g. for an integer or a memory address), or multiple bytes (e.g., for a string literal). The complete instruction set for CVM is given in the next section. Most instructions get their operands from the run-time stack. In general, the operands are removed from the stack whenever the instruction is executed, and any results are left on the top of the stack. With respect to boolean values, zero means false and any nonzero value is interpreted as true.

The following diagram illustrates a program loaded into memory.



Variable Addressing

Each time that a subprogram is called, CVM saves the current value of BP and sets BP to point to the new activation record (a.k.a., frame). When the subprogram returns, CVM restores BP back to the saved value.

A variable has an absolute address in memory (the address of the first byte), but variables are more commonly addressed relative to a register. A local variable is addressed relative to register BP, and a global variable is addressed relative to register SB.

Various load and store operations move data between memory and the run-time stack.

E.2 Implementation

CVM is implemented by three classes in package `edu.citadel.cvm`.

Class `Constants` defines the number of bytes for primitive types plus the number of bytes for the context part of an activation record.

```
public class Constants
{
    public static final int BYTES_PER_OPCODE = 1;
    public static final int BYTES_PER_INTEGER = 4;
    public static final int BYTES_PER_ADDRESS = 4;
    public static final int BYTES_PER_CHAR = 2;
    public static final int BYTES_PER_BOOLEAN = 1;
    public static final int BYTES_PER_CONTEXT = 2*BYTES_PER_ADDRESS;
}
```

Class Opcode is an enum class that defines the name and numeric values for each CVM opcode. In addition, this class defines several helper functions used by the CVM, the assembler, and disassembler, including static method toOpcode(byte b) that returns the opcode for the specified byte value.

```
public enum Opcode
{
    // halt opcode
    HALT(0),

    // load opcodes (move data from memory to top of stack)
    LOAD(9),
    LOADB(10),
    LOADW(12),
    LOADSTR(13),
    LDCB(14),
    LDCCH(15),
    LDCINT(16),
    ...

    // arithmetic opcodes
    ADD(70),
    SUB(71),
    MUL(72),
    DIV(73),
    ...

    // program/procedure opcodes
    PROGRAM(90),
    PROC(91),
    CALL(92),
    RET(93),
    ...

    private final byte value;

    /**
     * Construct an opcode with its machine instruction value.
     */
    private Opcode(int value)
    {
        this.value = (byte) value;
    }

    ...
}
```

Class CVM is the primary component of the implementation for the virtual machine. In addition to several helper methods, every opcode is implemented by a method. CVM method `run()` provides the basic control logic for the virtual machine using a large “switch” statement to dispatch opcodes to their corresponding method calls.

```
public void run()
{
    running = true;
    pc = 0;

    while (running)
    {
        switch (Opcode.toOpcode(fetchByte()))
        {
            case ADD      -> add();
            case ALLOC    -> allocate();
            case BR       -> branch();
            case BE       -> branchEqual();
            case BNE      -> branchNotEqual();
            case BG       -> branchGreater();
            case BGE      -> branchGreaterOrEqual();
            case BL       -> branchLess();
            case BLE      -> branchLessOrEqual();
            case BZ       -> branchZero();
            case BNZ      -> branchNonZero();
            case CALL     -> call();
            case DEC      -> decrement();
            case DIV      -> divide();
            case GETCH    -> getCh();
            case GETINT   -> getInt();
            case GETSTR   -> getString();
            case HALT     -> halt();
            case INC      -> increment();
            ...
            case SHL      -> shiftLeft();
            case SHR      -> shiftRight();
            case STORE    -> store();
            case STOREB   -> storeByte();
            case STORE2B  -> store2Bytes();
            case STOREW   -> storeWord();
            case SUB      -> subtract();
            default       -> error("invalid machine instruction");
        }
    }
}
```

E.3 CVM Instruction Set Architecture

Mnemonic	Short Description	Stack before <hr/> after	Definition
Arithmetic Opcodes			
ADD	Add: Pop two integers from the stack and push their sum back onto the stack.	$\begin{array}{r} n1 \\ n2 \\ \hline n1 + n2 \end{array}$	$n2 \leftarrow \text{popInt}()$ $n1 \leftarrow \text{popInt}()$ $\text{pushInt}(n1 + n2)$
SUB	Subtract: Pop two integers from the stack and push their difference back onto the stack.	$\begin{array}{r} n1 \\ n2 \\ \hline n1 - n2 \end{array}$	$n2 \leftarrow \text{popInt}()$ $n1 \leftarrow \text{popInt}()$ $\text{pushInt}(n1 - n2)$
MUL	Multiply: Pop two integers from the stack and push their product back onto the stack.	$\begin{array}{r} n1 \\ n2 \\ \hline n1 * n2 \end{array}$	$n2 \leftarrow \text{popInt}()$ $n1 \leftarrow \text{popInt}()$ $\text{pushInt}(n1 * n2)$
DIV	Divide: Pop two integers from the stack and push their quotient back onto the stack.	$\begin{array}{r} n1 \\ n2 \\ \hline n1 / n2 \end{array}$	$n2 \leftarrow \text{popInt}()$ $n1 \leftarrow \text{popInt}()$ $\text{pushInt}(n1 / n2)$
MOD	Modulo: Pop two integers from the stack, divide them, and push the remainder back onto the stack.	$\begin{array}{r} n1 \\ n2 \\ \hline n1 \% n2 \end{array}$	$n2 \leftarrow \text{popInt}()$ $n1 \leftarrow \text{popInt}()$ $\text{pushInt}(n1 \% n2)$
NEG	Negate: Pop an integer from the stack, negate it, and push the result back onto the stack.	$\begin{array}{r} n \\ \hline -n \end{array}$	$n \leftarrow \text{popInt}()$ $\text{pushInt}(-n)$
INC	Increment: Pop an integer from the stack, add 1, and push the result back onto the stack.	$\begin{array}{r} n \\ \hline n + 1 \end{array}$	$n \leftarrow \text{popInt}()$ $\text{pushInt}(n + 1)$
DEC	Decrement: Pop an integer from the stack, subtract 1, and push the result back onto the stack.	$\begin{array}{r} n \\ \hline n - 1 \end{array}$	$n \leftarrow \text{popInt}()$ $\text{pushInt}(n - 1)$

Logical Opcodes			
NOT	Logical Not: Pop a byte from the stack and push its logical negation back onto the stack.	$\frac{b}{!b}$	<pre> b ← popByte() if b = 0 pushByte(1) else pushByte(0) </pre>
Shift Opcodes			
SHL b	<p>Shift Left: Pop an integer from the stack, shift the bits left by the specified amount using zero fill, and push the result back onto the stack.</p> <p>Note: Only the right most five bits of the argument are used for the shift.</p>	$\frac{n}{n \ll b}$	<pre> n ← popInt() pushInt(n << b) </pre>
SHR b	<p>Shift Right: Pop an integer from the stack, shift it right by the specified amount using sign extend, and push the result back onto the stack.</p> <p>Note: Only the right most five bits of the argument are used for the shift.</p>	$\frac{n}{n \gg b}$	<pre> n ← popInt() pushInt(n >> b) </pre>
Branch Opcodes			
BR displ	Branch: Branch unconditionally according to displacement argument (may be positive or negative).		$pc \leftarrow pc + displ$

BE displ	Branch Equal: Pop two integers from the stack and compare them. If they are equal, then branch according to displacement argument (may be positive or negative); otherwise continue with the next instruction.	$\frac{n1}{n2}$	<pre> n2 ← popInt() n1 ← popInt() if n1 == n2 pc ← pc + displ </pre>
BNE displ	Branch Not Equal: Pop two integers from the stack and compare them. If they are not equal, then branch according to displacement argument (may be positive or negative); otherwise continue with the next instruction.	$\frac{n1}{n2}$	<pre> n2 ← popInt() n1 ← popInt() if n1 != n2 pc ← pc + displ </pre>
BG displ	Branch Greater: Pop two integers from the stack and compare them. If the second integer is greater than the first, then branch according to displacement argument (may be positive or negative); otherwise continue with the next instruction.	$\frac{n1}{n2}$	<pre> n2 ← popInt() n1 ← popInt() if n1 > n2 pc ← pc + displ </pre>
BGE displ	Branch Greater or Equal: Pop two integers from the stack and compare them. If the second integer is greater than or equal to the first, then branch according to displacement argument (may be positive or negative); otherwise continue with the next instruction.	$\frac{n1}{n2}$	<pre> n2 ← popInt() n1 ← popInt() if n1 >= n2 pc ← pc + displ </pre>

BL displ	Branch Less: Pop two integers from the stack and compare them. If the second integer is less than the first, then branch according to displacement argument (may be positive or negative); otherwise continue with the next instruction.	$\frac{n1}{n2}$	<pre> n2 ← popInt() n1 ← popInt() if n1 < n2 pc ← pc + displ </pre>
BLE displ	Branch Less or Equal: Pop two integers from the stack and compare them. If the second integer is less than or equal to the first, then branch according to displacement argument (may be positive or negative); otherwise continue with the next instruction.	$\frac{n1}{n2}$	<pre> n2 ← popInt() n1 ← popInt() if n1 ≤ n2 pc ← pc + displ </pre>
BZ displ	Branch if Zero: Pop one byte from the stack. If it is zero, then branch according to displacement argument (may be positive or negative); otherwise continue with the next instruction.	$\frac{b}{}$	<pre> b ← popByte() if b = 0 pc ← pc + displ </pre>
BNZ displ	Branch if Nonzero: Pop one byte from the stack. If it is nonzero then branch according to displacement argument (may be positive or negative); otherwise continue with the next instruction.	$\frac{b}{}$	<pre> b ← popByte() if b ≠ 0 pc ← pc + displ </pre>
Load/Store Opcodes			
LOAD n	Load multiple bytes onto the stack: The number of bytes to move is part of the instruction. Pop an address from the stack and push n bytes starting at that address onto the stack.	$\frac{\text{addr}}{ \begin{array}{l} b1 \\ b2 \\ \cdots \\ bn \end{array} }$	<pre> addr ← popInt(); for i ← 0..n-1 loop pushByte(mem[addr + i]) </pre>

LOADB	Load Byte: Load (push) a single byte onto the stack. The address of the byte is obtained by popping it off the stack.	$\frac{\text{addr}}{b}$	$\begin{aligned} \text{addr} &\leftarrow \text{popInt}() \\ b &\leftarrow \text{mem}[\text{addr}] \\ \text{pushByte}(b) \end{aligned}$
LOAD2B	Load Two Bytes: Load (push) two consecutive bytes onto the stack. The address of the first byte is obtained by popping it off the stack.	$\frac{\text{addr}}{\begin{matrix} b0 \\ b1 \end{matrix}}$	$\begin{aligned} \text{addr} &\leftarrow \text{popInt}() \\ b0 &\leftarrow \text{mem}[\text{addr} + 0] \\ b1 &\leftarrow \text{mem}[\text{addr} + 1] \\ \text{pushByte}(b0) \\ \text{pushByte}(b1) \end{aligned}$
LOADW	Load Word: Load (push) a word (four consecutive bytes) onto the stack. The address of the word is obtained by popping it off the stack.	$\frac{\text{addr}}{w}$	$\begin{aligned} \text{addr} &\leftarrow \text{popInt}() \\ w &\leftarrow \text{getWord}(\text{addr}) \\ \text{pushInt}(w) \end{aligned}$
LDCB <i>b</i>	Load Constant Byte: Fetch the byte immediately following the opcode and push it onto the stack.	$\frac{\quad}{b}$	$\text{pushByte}(b)$
LDCB0	Load Constant Byte 0: Optimized version of LDCB 0.	$\frac{\quad}{0}$	$\text{pushByte}(0)$
LDCB1	Load Constant Byte 1: Optimized version of LDCB 1.	$\frac{\quad}{1}$	$\text{pushByte}(1)$
LDCCH <i>c</i>	Load Constant Character: Fetch the character immediately following the opcode and push it onto the stack.	$\frac{\quad}{c}$	$\text{pushChar}(c)$
LDCINT <i>n</i>	Load Constant Integer: Fetch the integer immediately following the opcode and push it onto the stack.	$\frac{\quad}{n}$	$\text{pushInt}(n)$
LDCINT0	Load Constant Integer 0: Optimized version of LDCINT 0.	$\frac{\quad}{0}$	$\text{pushInt}(0)$
LDCINT1	Load Constant Integer 1: Optimized version of LDCINT 1.	$\frac{\quad}{1}$	$\text{pushInt}(1)$

LDCSTR s	Load Constant String: The string (length plus characters) immediately follows the opcode. Push the string (length and characters) onto the stack.	$\frac{s}{s}$	$n \leftarrow \text{fetchInt}()$ $\text{pushInt}(n)$ for $i \leftarrow 0..n-1$ loop $c \leftarrow \text{fetchChar}()$ $\text{pushChar}(c)$
LDLADDR n	Load Local Address: Compute the absolute address of a local variable from its relative address n and push the absolute address onto the stack.	$\frac{bp + n}{bp + n}$	$\text{pushInt}(bp + n)$
LDGADDR n	Load Global Address: Compute the absolute address of a global (program level) variable from its relative address n and push the absolute address onto the stack.	$\frac{sb + n}{sb + n}$	$\text{pushInt}(sb + n)$
STORE n	Store n bytes: Remove n bytes from the stack followed by an absolute address and copy the n bytes to the location starting at the absolute address.	$\frac{\begin{matrix} \text{addr} \\ b1 \\ b2 \\ \dots \\ bn \end{matrix}}{\text{addr} \\ b1 \\ b2 \\ \dots \\ bn}$	for $i \leftarrow n-1..0$ loop $\text{data}[i] \leftarrow \text{popByte}()$ $\text{addr} \leftarrow \text{popInt}()$ for $i \leftarrow 0..n-1$ loop $\text{mem}[\text{addr} + i] \leftarrow \text{data}[i]$
STOREB	Store Byte: Store a single byte at a specified memory location. The byte to be stored and the address where it is to be stored are popped from the stack.	$\frac{\text{addr} \quad b}{\text{addr} \quad b}$	$b \leftarrow \text{popByte}()$ $\text{addr} \leftarrow \text{popInt}()$ $\text{memory}[\text{addr}] \leftarrow b$
STORE2B	Store Two Bytes: Store two bytes at a specified memory location. The bytes to be stored and the address where they are to be stored are popped from the stack.	$\frac{\text{addr} \quad b0 \quad b1}{\text{addr} \quad b0 \quad b1}$	$b1 \leftarrow \text{popByte}()$ $b0 \leftarrow \text{popByte}()$ $\text{addr} \leftarrow \text{popInt}()$ $\text{mem}[\text{addr} + 0] \leftarrow b0$ $\text{mem}[\text{addr} + 1] \leftarrow b1$
STOREW	Store Word: Store a word (4 bytes) at a specified memory location. The word to be stored and the address where it is to be stored are popped from the stack.	$\frac{\text{addr} \quad w}{\text{addr} \quad w}$	$w \leftarrow \text{popWord}()$ $\text{addr} \leftarrow \text{popInt}()$ $\text{putWord}(w, \text{addr})$

ALLOC n	Allocate: Allocate space on the stack for future use.		$sp \leftarrow sp + n$
Program/Procedure Opcodes			
PROGRAM n	Program: Initialize base pointer and allocate space on the stack for the program's local variables.		$bp \leftarrow sb$ $sp \leftarrow bp + n - 1$
PROC n	Procedure: Allocate space on the stack for a subprogram's local variables.		$sp \leftarrow sp + n$
CALL $disp$	Call: Call a subprogram, pushing current values for BP and PC onto the stack.	$\frac{bp}{p}$	$pushInt(bp)$ $pushInt(pc)$ $bp \leftarrow sp - 7$ $pc \leftarrow pc + disp$
RET n	Return: Return from a subprogram, restoring the old value for BP plus space on stack previously allocated for the subprogram's local variables and parameters.		$bpSave \leftarrow bp$ $sp \leftarrow bpSave - n - 1$ $bp \leftarrow getInt(bpSave)$ $pc \leftarrow getInt(bpSave + 4)$
RET0	Optimized version of RET 0.		$bpSave \leftarrow bp$ $sp \leftarrow bpSave - 1$ $bp \leftarrow getInt(bpSave)$ $pc \leftarrow getInt(bpSave + 4)$
RET4	Optimized version of RET 4.		$bpSave \leftarrow bp$ $sp \leftarrow bpSave - 5$ $bp \leftarrow getInt(bpSave)$ $pc \leftarrow getInt(bpSave + 4)$
HALT	Halt: Stop the virtual machine.		halt
I/O Opcodes			
GETINT	Get Integer: Read digits from standard input, convert them to an integer, and store the integer at the address on top of stack.	$\frac{addr}{}$	$addr \leftarrow popInt()$ $n \leftarrow readInt()$ $putInt(n, addr)$

GETCH	Get Character: Read character from standard input and store it at the address on top of stack.	<u>addr</u>	$addr \leftarrow \text{popInt}()$ $c \leftarrow \text{readChar}()$ $\text{putChar}(c, \text{addr})$
GETSTR n	Get String: Read string from standard input and store it at the address on top of stack.	<u>addr</u>	$addr \leftarrow \text{popInt}()$ $s \leftarrow \text{readStr}()$ $\text{strLen} = \min(s.\text{length}, n)$ $\text{putInt}(n, \text{addr})$ for $i \leftarrow 0..\text{strLen}-1$ loop $\text{putChar}(s[i], \text{addr})$ $addr \leftarrow \text{addr} + 2$
PUTBYTE	Put Byte: Pop byte from top of the stack and write its value to standard output.	<u>b</u>	$b \leftarrow \text{popByte}()$ $\text{writeByte}(b)$
PUTINT	Put Integer: Pop integer from top of the stack and write its value to standard output.	<u>n</u>	$n \leftarrow \text{popInt}()$ $\text{writeInt}(n)$
PUTCH	Put Character: Pop character from top of stack and write its value to standard output.	<u>c</u>	$c \leftarrow \text{popChar}()$ $\text{writeChar}(c)$
PUTSTR n	Put String: Write a string of n characters to standard output. The string (length plus characters) was previously pushed onto the stack.	<u>s</u>	$n\text{Bytes} \leftarrow 4 + 2*n$ $addr \leftarrow sp - n\text{Bytes} + 1$ $\text{strLen} \leftarrow \text{getInt}(addr)$ for $i \leftarrow 0..\text{strLen}-1$ loop $\text{writeChar}(\text{mem}[2*i])$ $sp \leftarrow sp - n$
PUTEOL	Put End-of-Line: Write a line terminator to standard output.		$\text{write}(\text{EOL})$