

**Script** generated by TTT

Title: Petter: Compilerbau (27.07.2017)

Date: Thu Jul 27 14:15:17 CEST 2017

Duration: 97:31 min

Pages: 41

Topic:

Code Synthesis

Code Synthesis

## Chapter 1: The Register C-Machine

### The Register C-Machine (R-CMa)

We generate Code for the Register C-Machine.

The Register C-Machine is a virtual machine (VM).

- there exists no processor that can execute its instructions
- ... but we can build an interpreter for it
- we provide a visualization environment for the R-CMa
- the R-CMa has no **double**, **float**, **char**, **short** or **long** types
- the R-CMa has no instructions to communicate with the operating system
- the R-CMa has an unlimited supply of registers

## Virtual Machines

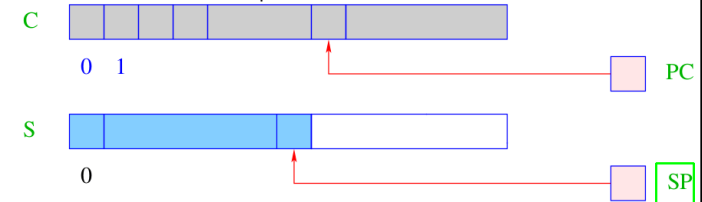
A virtual machine has the following ingredients:

- any virtual machine provides a set of **instructions**
- instructions are executed on virtual hardware
- the virtual hardware is a collection of **data structures** that is accessed and modified by the VM instructions
- ... and also by other components of the **run-time system**, namely functions that go beyond the instruction semantics
- the **interpreter** is part of the run-time system

245 / 289

## Components of a Virtual Machine

Consider **Java** as an example:



A virtual machine such as the **Dalvik VM** has the following structure:

- **S**: the data store – a memory region in which cells can be stored in LIFO order  $\rightsquigarrow$  **stack**.
- **SP**: ( $\hat{=}$  **stack pointer**) pointer to the last used cell in **S**
- beyond **S** follows the memory containing the heap

246 / 289

## Executing a Program

- the machine loads an instruction from **C[PC]** into the **instruction register IR** in order to execute it
- before evaluating the instruction, the **PC** is incremented by one

```
while (true) {  
  IR = C[PC]; PC++;  
  execute (IR);  
}
```

- note: the **PC** must be incremented **before** the execution, since an instruction may modify the **PC**
- the loop is exited by evaluating a **halt** instruction that returns directly to the operating system

247 / 289

Code Synthesis



Chapter 2:

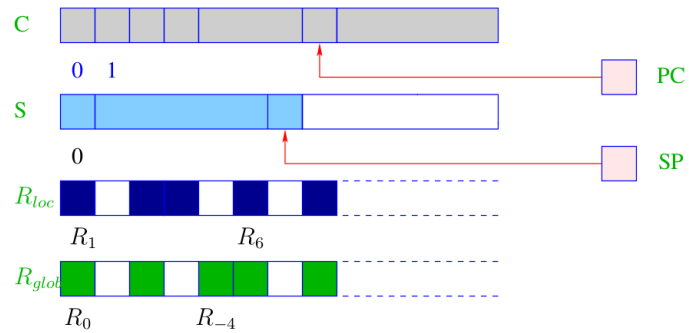
Generating Code for the Register C-Machine

248 / 289

## Principles of the R-CMa

The R-CMa is composed of a stack, heap and a code segment, just like the JVM; it additionally has register sets:

- *local* registers are  $R_1, R_2, \dots, R_i, \dots$
- *global* registers are  $R_0, R_{-1}, \dots, R_{-n}, \dots$



250 / 289

## The Register Sets of the R-CMa

The two register sets have the following purpose:

- the *local* registers  $R_i$ 
  - save temporary results
  - store the contents of local variables of a function
  - can efficiently be stored and restored from the stack

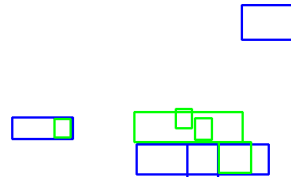


251 / 289

## Translation of Simple Expressions

Using variables stored in registers; loading constants:

instruction	semantics	intuition
load $R_i$ $c$	$R_i = c$	load constant
move $R_i$ $R_j$	$R_i = R_j$	copy $R_j$ to $R_i$



252 / 289

## Translation of Expressions

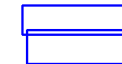
Let  $op = \{add, sub, div, mul, mod, le, gr, eq, leq, geq, and, or\}$ .  
The R-CMa provides an instruction for each operator  $op$ .

$$op \ R_i \ R_j \ R_k$$

where  $R_i$  is the target register,  $R_j$  the first and  $R_k$  the second argument.

Correspondingly, we generate code as follows:

$$code_R^i \ e_1 \ op \ e_2 \ \rho = \begin{matrix} code_R^i \ e_1 \ \rho \\ code_R^{i+1} \ e_2 \ \rho \\ op \ R_i \ R_i \ R_{i+1} \end{matrix}$$



253 / 289

## Managing Temporary Registers

Observe that temporary registers are re-used: translate  $3*4+3*4$  with  $t = 4$ :

$$\text{code}_R^4 \ 3*4+3*4 \ \rho = \begin{array}{l} \text{code}_R^4 \ 3*4 \ \rho \\ \text{code}_R^5 \ 3*4 \ \rho \\ \text{add } R_4 \ R_4 \ R_5 \end{array}$$

where

$$\text{code}_R^i \ 3*4 \ \rho = \begin{array}{l} \text{loadc } R_i \ 3 \\ \text{loadc } R_{i+1} \ 4 \\ \text{mul } R_i \ R_i \ R_{i+1} \end{array}$$

we obtain

$$\text{code}_R^4 \ 3*4+3*4 \ \rho =$$


254 / 289

## Semantics of Operators

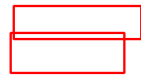
The operators have the following semantics:

<code>add</code> $R_i \ R_j \ R_k$	$R_i = R_j + R_k$
<code>sub</code> $R_i \ R_j \ R_k$	$R_i = R_j - R_k$
<code>div</code> $R_i \ R_j \ R_k$	$R_i = R_j / R_k$
<code>mul</code> $R_i \ R_j \ R_k$	$R_i = R_j * R_k$
<code>mod</code> $R_i \ R_j \ R_k$	$R_i = \text{signum}(R_k) \cdot k$ with $ R_j  = n \cdot  R_k  + k \wedge n \geq 0, 0 \leq k <  R_k $
<code>le</code> $R_i \ R_j \ R_k$	$R_i = \text{if } R_j < R_k \text{ then } 1 \text{ else } 0$
<code>gr</code> $R_i \ R_j \ R_k$	$R_i = \text{if } R_j > R_k \text{ then } 1 \text{ else } 0$
<code>eq</code> $R_i \ R_j \ R_k$	$R_i = \text{if } R_j = R_k \text{ then } 1 \text{ else } 0$
<code>leq</code> $R_i \ R_j \ R_k$	$R_i = \text{if } R_j \leq R_k \text{ then } 1 \text{ else } 0$
<code>geq</code> $R_i \ R_j \ R_k$	$R_i = \text{if } R_j \geq R_k \text{ then } 1 \text{ else } 0$
<code>and</code> $R_i \ R_j \ R_k$	$R_i = R_j \ \& \ R_k$ // bit-wise and
<code>or</code> $R_i \ R_j \ R_k$	$R_i = R_j \   \ R_k$ // bit-wise or

255 / 289

## Translation of Unary Operators

Unary operators  $\text{op} = \{\text{neg}, \text{not}\}$  take only two registers:

$$\text{code}_R^i \ \text{op } e \ \rho = \begin{array}{l} \text{code}_R^i \ e \ \rho \\ \text{op } R_i \ R_i \end{array}$$


256 / 289

## Applying Translation Schema for Expressions

Suppose the following function `void f(void)` is given:

```
void f(void) {
    int x, y, z;
    x = y+z*3;
}
```

- Let  $\rho = \{x \mapsto 1, y \mapsto 2, z \mapsto 3\}$  be the address environment.
- Let  $R_4$  be the first free register, that is,  $i = 4$ .

$$\text{code}_R^4 \ x=y+z*3 \ \rho = \begin{array}{l} \text{code}_R^4 \ y+z*3 \ \rho \\ \text{move } R_1 \ R_4 \end{array}$$

257 / 289

### Chapter 3: Statements and Control Structures



### Translation of Statement Sequences

The code for a sequence of statements is the concatenation of the instructions for each statement in that sequence:

$$\begin{aligned}
 \text{code}^i (s \text{ } ss) \rho &= \text{code}^i s \rho \quad \text{code}^i ss \rho \\
 \text{code}^i \varepsilon \rho &= \quad \quad \quad // \text{ empty sequence of instructions}
 \end{aligned}$$

Note here:  $s$  is a statement,  $ss$  is a sequence of statements

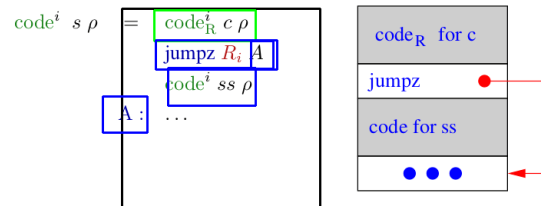


### Simple Conditional

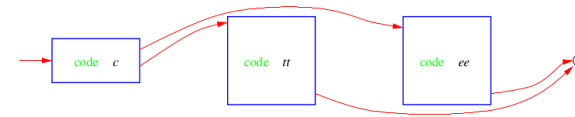
We first consider  $s \equiv \text{if}(c) \text{ } ss.$  ...and present a translation without basic blocks.

Idea:

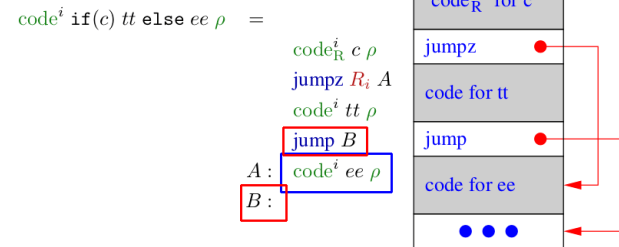
- emit the code of  $c$  and  $ss$  in sequence
- insert a jump instruction in-between, so that correct control flow is ensured



### General Conditional



Translation of  $\text{if}(c) \text{ } tt \text{ } \text{else } ee.$



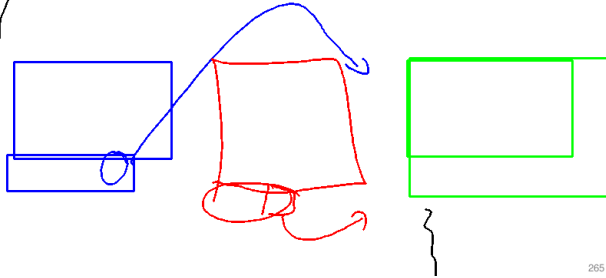
## Example for if-statement

Let  $\rho = \{x \mapsto 4, y \mapsto 7\}$  and let  $s$  be the statement

```

if (x > y) { /* (i) */
  x = x - y; /* (ii) */
} else { /* (iii) */
  y = y - x;
}
    
```

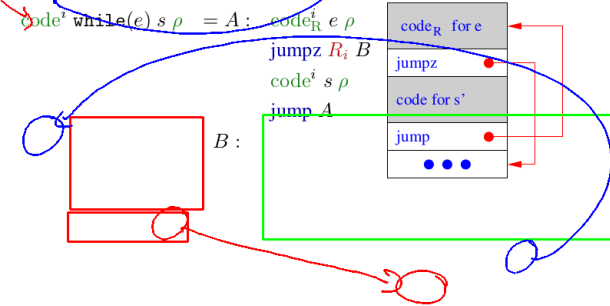
Then  $\text{code}^i s \rho$  yields:



265 / 289

## Iterating Statements

We only consider the loop  $s \equiv \text{while}(e) s'$ . For this statement we define:



266 / 289

## for-Loops

The for-loop  $s \equiv \text{for}(e_1; e_2; e_3) s'$  is equivalent to the statement sequence  $e_1; \text{while}(e_2) \{s'; e_3;\}$  – as long as  $s'$  does not contain a **continue** statement.

Thus, we translate:

```

code^i for(e1; e2; e3) s rho = code_R^i e1 rho
                             A: code_R^i e2 rho
                               jumpz R_i B
                               code^i s rho
                               code_R^i e3 rho
                               jump A
                             B:
    
```

268 / 289

## Consecutive Alternatives

Let **switch**  $s$  be given with  $k$  consecutive **case** alternatives:

```

switch (e) {
  case 0: s_0; break;
  :
  case k-1: s_{k-1}; break;
  default: s_k; break;
}
    
```



270 / 289

## Translation of the $check^i$ Macro

The macro  $check^i l u B$  checks if  $l \leq R_i < u$ . Let  $k = u - l$ .

- if  $l \leq R_i < u$  it jumps to  $B + R_i - l$
- if  $R_i < l$  or  $R_i \geq u$  it jumps to  $A_k$



```

B:  jump A_0
    :
    :
    :      jump A_k
C:

```

271 / 289

## Translation of the $check^i$ Macro

The macro  $check^i l u B$  checks if  $l \leq R_i < u$ . Let  $k = u - l$ .

- if  $l \leq R_i < u$  it jumps to  $B + R_i - l$
- if  $R_i < l$  or  $R_i \geq u$  it jumps to  $A_k$

we define:

```

check^i l u B =  loadc R_{i+1} l
                  geq R_{i+2} R_i R_{i+1}
                  jumpz R_{i+2} E      B:  jump A_0
                  sub  R_i R_i R_{i+1}
                  loadc R_{i+1} u
                  geq R_{i+2} R_i R_{i+1}
                  jumpz R_{i+2} D      :      :
                  :
                  :      jump A_k
E:  loadc R_i u - l
D:  jumpi R_i B

```

271 / 289

## Improvements for Jump Tables

This translation is only suitable for *certain* switch-statement.

- In case the table starts with 0 instead of  $u$  we don't need to subtract it from  $e$  before we use it as index
- if the value of  $e$  is **guaranteed** to be in the interval  $[l, u]$ , we can omit  $check$

272 / 289

## General translation of switch-Statements

In general, the values of the various cases may be far apart:

- generate an **if**-ladder, that is, a sequence of **if**-statements
- for  $n$  cases, an **if**-cascade (tree of conditionals) can be generated  $\sim O(\log n)$  tests
- if the sequence of numbers has small gaps ( $\leq 3$ ), a jump table may be smaller and faster
- one could generate several jump tables, one for each sets of consecutive cases
- an **if** cascade can be re-arranged by using information from *profiling*, so that paths executed more frequently require fewer tests

273 / 289

## Chapter 4: Functions

274 / 289

## Ingredients of a Function

The definition of a function consists of

- a **name** with which it can be called;
- a specification of its **formal parameters**;
- possibly a **result type**;
- a sequence of **statements**.

In C we have:

$$\text{code}_R^i f \rho = \text{loadc } R_i \_f \text{ with } \_f \text{ starting address of } f$$

Observe:

- function names must have an address assigned to them
- since the size of functions is unknown before they are translated, the addresses of forward-declared functions must be inserted later

275 / 289

## Memory Management in Functions

```

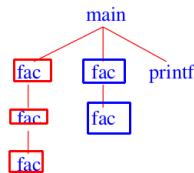
int fac(int x) {
    if (x<=0) return 1;
    else return x*fac(x-1);
}

int main(void) {
    int n;
    n = fac(2) + fac(1);
    printf("%d", n);
}

```

At run-time several **instances** may be active, that is, the function has been called but has not yet returned.

The recursion tree in the example:



276 / 289

## Memory Management in Function Variables

The **formal parameters** and the **local variables** of the various **instances** of a function must be kept separate

**Idea for implementing functions:**

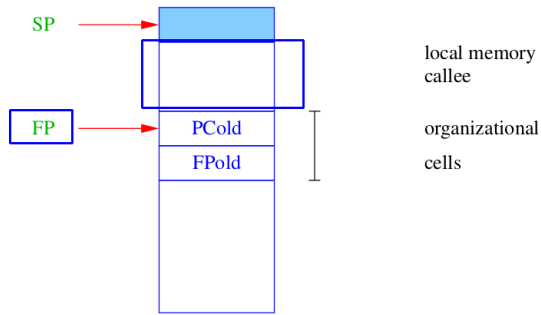
- set up a region of memory each time it is called
- in sequential programs this memory region can be allocated on the stack
- thus, each instance of a **function has its own region on the stack**
- these regions are called **stack frames**

277 / 289



## Organization of a Stack Frame

- stack representation: grows upwards
- SP points to the last used stack cell



278 / 289

## Split of Obligations

### Definition

Let  $f$  be the current function that calls a function  $g$ .

- $f$  is dubbed *caller*
- $g$  is dubbed *callee*

The code for managing function calls has to be split between caller and callee.

This split cannot be done arbitrarily since some information is only known in that caller or only in the callee.

### Observation:

The space requirement for parameters is only known by the caller:

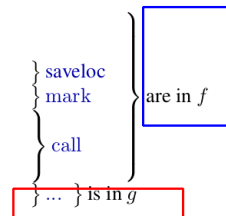
Example: `printf`

279 / 289

## Principle of Function Call and Return

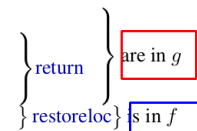
actions taken on entering  $g$ :

1. compute the start address of  $g$
2. compute actual parameters in globals
3. backup of caller-save registers
4. backup of FP
5. set the new FP
6. back up of PC and jump to the beginning of  $g$
7. copy actual params to locals



actions taken on leaving  $g$ :

1. compute the result into  $R_0$
2. restore FP, SP
3. return to the call site in  $f$ , that is, restore PC
4. restore the caller-save registers



280 / 289

## Managing Registers during Function Calls

The two register sets (global and local) are used as follows:

- automatic variables live in *local* registers  $R_i$
- intermediate results also live in *local* registers  $R_i$
- parameters live in *global* registers  $R_i$  (with  $i \leq 0$ )
- global variables:

281 / 289

## Translation of Function Calls

A function call  $g(e_1, \dots, e_n)$  is translated as follows:

```

codeRi g(e1, ..., en) ρ = codeRi g ρ
                          codeRi+1 e1 ρ
                          ⋮
                          codeRi+n en ρ
                          move R-1 Ri+1
                          ⋮
                          move R-n Ri+n
                          saveloc R1 Ri-1
                          mark
                          call Ri
                          restoreloc R1 Ri-1
                          move Ri R0
    
```

282 / 289

## Rescuing the FP

The instruction `mark` allocates stack space for the return value and the organizational cells and backs up FP.



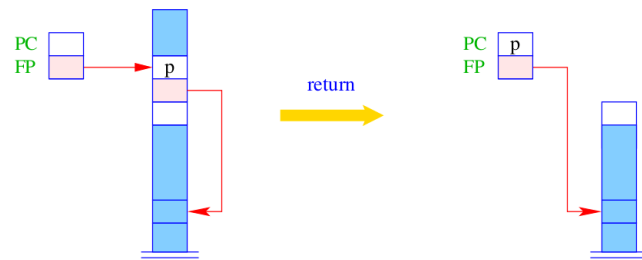
```

S[SP+1] = FP;
SP = SP + 1;
    
```

283 / 289

## Return from a Function

The instruction `return` relinquishes control of the current stack frame, that is, it restores PC and FP.



```

PC = S[FP];
SP = FP-2;
FP = S[SP+1];
    
```

286 / 289

## Translation of Whole Programs

A program  $P = F_1; \dots; F_n$  must have a single `main` function.

```

code1 P ρ = load R1 _main
             mark
             call R1
             halt
             _f1 : code1 F1 ρ ⊕ ρf1
             ⋮
             _fn : code1 Fn ρ ⊕ ρfn
    
```

288 / 289

## Translation of the fac-function

Consider:

```
int fac(int x) {
  if (x<=0) then
    return 1;
  else
    return x*fac(x-1);
}

_fac:  move R1 R_{-1}  save param.
i = 2  move R2 R1     if (x<=0)
      load R3 0
      leq R2 R2 R3
      jumpz R2 _A   to else
      load R2 1    return 1
      move R0 R2
      return
      jump _B      code is dead

_A:    move R2 R1     x*fac(x-1)
i = 3  move R3 R1     x-1
i = 4  load R4 1
      sub R3 R3 R4
      move R_{-1} R3  fac(x-1)
i = 3  load R3 _fac
      saveloc R1 R2
      mark
      call R3
      restoreloc R1 R2
      move R3 R0
      mul R2 R2 R3
      move R0 R2     return x*...
      return
      return
_B:    return
```

289 / 289

Ende der Präsentation. Klicken Sie zum Schließen.