

G53CMP: Lecture 12 & 13

Code Generation II

Henrik Nilsson

University of Nottingham, UK

G53CMP: Lecture 12 & 13 – p.1/46

Code Generation: Demo II (1)

And for this program using arrays and a procedure:

```
let
  proc swap(var x: Integer, var y: Integer)
    let
      var t: Integer
    in begin
      t := x; x := y; y := t
    end;
  var a: Integer[5] := [7,3,1,9,2];
  var i: Integer;
  var j: Integer
```

G53CMP: Lecture 12 & 13 – p.3/46

Code Generation: Demo I

Let us generate code for:

```
let
  var f: Integer := 1;
  var i: Integer := 1
in
  while i <= 10 do begin
    f := f * i;
    putint(f);
    i := i + 1
  end
```

G53CMP: Lecture 12 & 13 – p.2/46

Code Generation: Demo II (2)

```
in begin
  i := 0;
  while i < 4 do begin
    j := i + 1;
    while j < 5 do begin
      if a[i] > a[j] then
        swap(a[i], a[j])
      else skip();
      j := j + 1
    end;
    i := i + 1
  end;
```

G53CMP: Lecture 12 & 13 – p.4/46

Code Generation: Demo II (3)

```

i := 0;
while i <= 4 do begin
  putint(a[i]);
  i := i + 1
end
end

```

G53CMP: Lecture 12 & 13 – p.5/46

Specifying Code Selection (2)

Note:

- *execute* **generates code** for executing a command (it does not execute a command directly);
- *evaluate* **generates code** for evaluating an expression, leaving the result on the top of the stack.
- *elaborate* **generates code** for reserving storage for declared constants and variables, evaluating any initialisation expressions, and for declared procedures and functions.

G53CMP: Lecture 12 & 13 – p.7/46

Specifying Code Selection (1)

- Code selection is specified **inductively** over the phrases of the source language:

$$\begin{array}{l}
 \text{Command} \rightarrow \text{Identifier} := \text{Expression} \\
 \quad \quad \quad | \text{Command} ; \text{Command} \\
 \dots
 \end{array}$$

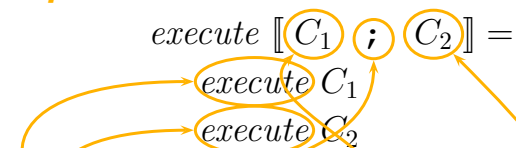
- **Code Function**: maps a source phrase to an instruction sequence. For example:

$$\begin{array}{l}
 \text{execute} : \text{Command} \rightarrow \text{Instruction}^* \\
 \text{evaluate} : \text{Expression} \rightarrow \text{Instruction}^* \\
 \text{elaborate} : \text{Declaration} \rightarrow \text{Instruction}^*
 \end{array}$$

G53CMP: Lecture 12 & 13 – p.6/46

Specifying Code Selection (3)

- Code functions are specified by means of **code templates**:



- The brackets \llbracket and \rrbracket enclose pieces of **concrete syntax** and **meta variables**.
- Note the **recursion**; i.e. inductive definition over the underlying phrase structure.

(Think of $\llbracket \cdot \rrbracket$ as a map from concrete to abstract syntax as specified by the abstract syntax grammars.)

G53CMP: Lecture 12 & 13 – p.8/46

Specifying Code Selection (4)

In a simple language, the code template for assignment might be:

$$\begin{aligned} \text{execute } \llbracket I := E \rrbracket = & \\ & \text{evaluate } E \\ & \text{STORE } \text{addr}(I) \end{aligned}$$

where

$$\text{addr} : \text{Identifier} \rightarrow \text{Address}$$

Note that the instruction sequences and individual instructions in the RHS of the defining equation are implicitly **concatenated**.

Note: meta variables range over **abstract syntax**.

G53CMP: Lecture 12 & 13 – p.9/46

Not Quite that Simple ...

However, something is clearly missing! Recall:

$$\begin{aligned} \text{execute} & : \text{Command} \rightarrow \text{Instruction}^* \\ \text{evaluate} & : \text{Expression} \rightarrow \text{Instruction}^* \\ \text{elaborate} & : \text{Declaration} \rightarrow \text{Instruction}^* \\ \text{addr} & : \text{Identifier} \rightarrow \text{Address} \end{aligned}$$

and consider again:

$$\begin{aligned} \text{execute } \llbracket I := E \rrbracket = & \\ & \text{evaluate } E \\ & \text{STORE } \text{addr}(I) \end{aligned}$$

How can the function *addr* possibly map an identifier (a name) to an address?

G53CMP: Lecture 12 & 13 – p.11/46

Exercise: Code Templates

Generate code for the fragment

$$\begin{aligned} f & := f * n; \\ n & := n - 1 \end{aligned}$$

using the following two templates:

$$\begin{aligned} \text{execute } \llbracket C_1 ; C_2 \rrbracket = & \quad \text{execute } \llbracket I := E \rrbracket = \\ & \text{execute } C_1 \quad \quad \quad \text{evaluate } E \\ & \text{execute } C_2 \quad \quad \quad \text{STORE } \text{addr}(I) \end{aligned}$$

and $\text{addr}(f) = [\text{SB} + 11]$, $\text{addr}(n) = [\text{SB} + 17]$.

Expand as far as the above templates allow.

G53CMP: Lecture 12 & 13 – p.10/46

Not Quite that Simple ... (2)

In more detail:

- *elaborate* is responsible for **assigning** addresses to variables
- a function like *addr* needs **access** to the addresses assigned by *elaborate*
- but the given type signatures for the code functions do **not permit** this communication!

G53CMP: Lecture 12 & 13 – p.12/46

Not Quite that Simple ... (3)

Consequently:

- The code functions need an additional **stack environment argument**, associating variables with addresses.
- The code function *elaborate* must **return an updated stack environment**.
- Need to keep track of the **current stack depth** (with respect to LB) to allow *elaborate* to determine the address (within activation record) for a new variable.

G53CMP: Lecture 12 & 13 – p.13/46

Not Quite that Simple ... (6)

To clearly convey the basic ideas first, we will:

- Use simplified MiniTriangle as main example:
 - No user-defined procedures or functions (only predefined, global ones).
 - Consequently, all variables are global (addressed with respect to SB).
 - No arrays (only simple variables, all of size 1 word) .
- Gloss over the bookkeeping details for the most part.

G53CMP: Lecture 12 & 13 – p.15/46

Not Quite that Simple ... (5)

- Need to keep track of the current scope level as the difference of current scope level and the scope level of a variable is needed in addition to its address to access it (see later lecture on run-time organisation and **static links**).

Moreover, need to generate **fresh names** for jump targets (recall the demo).

G53CMP: Lecture 12 & 13 – p.14/46

Not Quite that Simple ... (7)

However:

- Additional details will be given occasionally.
- Will revisit at appropriate points in lectures on run-time organisation.
- Refer to the HMTCC (coursework compiler) source code for full details.

G53CMP: Lecture 12 & 13 – p.16/46

Code Functions for MiniTriangle

In the simplified exposition, we can consider the code functions to have the following types:

```
run : Program → Instruction*
execute : Command → Instruction*
execute* : Command* → Instruction*
evaluate : Expression → Instruction*
evaluate* : Expression* → Instruction*
fetch : Identifier → Instruction*
assign : Identifier → Instruction*
elaborate : Declaration → Instruction*
elaborate* : Declaration* → Instruction*
```

G53CMP: Lecture 12 & 13 – p.17/46

Some HMTC Code Functions

```
execute :: Level -> CGEnv -> Depth -> Command
        -> CG TAMInst ()
```

```
evaluate :: Level -> CGEnv -> Expression
        -> CG TAMInst ()
```

```
elaborateDecls :: Level -> CGEnv -> Depth
               -> [Declaration]
               -> CG TAMInst (CGEnv, Depth)
```

(In essence: actual signatures differ in minor ways.)

G53CMP: Lecture 12 & 13 – p.19/46

A Code Generation Monad

HMTC uses a *Code Generation monad* to facilitate some of the bookkeeping:

```
instance Monad (CG instr)
```

Takes care of:

- Collation of generated instructions
- Generation of fresh names

Typical operations:

- `emit :: instr -> CG instr ()`
- `newName :: CG instr Name`

G53CMP: Lecture 12 & 13 – p.18/46

MiniTriangle Abstract Syntax Part I

(Simplified: no procedures, functions, arrays)

<i>Program</i>	→	<i>Command</i>	<i>Program</i>
<i>Command</i>	→	<i>Identifier</i> := <i>Expression</i>	CmdAssign
		<i>Identifier</i> (<i>Expression</i> *)	CmdCall
		begin <i>Command</i> * end	CmdSeq
		if <i>Expression</i> then <i>Command</i>	CmdIf
		else <i>Command</i>	
		while <i>Expression</i> do <i>Command</i>	CmdWhile
		let <i>Declaration</i> * in <i>Command</i>	CmdLet

G53CMP: Lecture 12 & 13 – p.20/46

Meta Variable Conventions

$C \in \text{Command}$
 $Cs \in \text{Command}^*$
 $E \in \text{Expression}$
 $Es \in \text{Expression}^*$
 $D \in \text{Declaration}$
 $Ds \in \text{Declaration}^*$
 $I \in \text{Identifier}$
 $O \in \text{Operator}$
 $IL \in \text{IntegerLiteral}$
 $TD \in \text{TypeDenoter}$

G53CMP: Lecture 12 & 13 – p.21/46

Code Function *execute* (1)

$run : \text{Program} \rightarrow \text{Instruction}^*$
 $execute : \text{Command} \rightarrow \text{Instruction}^*$

$run \llbracket C \rrbracket =$
 $execute \ C$
HALT

$execute \llbracket I := E \rrbracket =$
 $evaluate \ E$
 $assign \ I$

G53CMP: Lecture 12 & 13 – p.22/46

Code Function *execute* (2)

$execute \llbracket I (Es) \rrbracket =$
 $evaluate^* \ Es$
CALL $addr(I)$

$execute \llbracket begin \ Cs \ end \rrbracket =$
 $execute^* \ Cs$

G53CMP: Lecture 12 & 13 – p.23/46

Code Function *execute* (3)

$execute \llbracket \text{if } E \text{ then } C_1 \text{ else } C_2 \rrbracket =$
 $evaluate \ E$
JUMPIFZ g
 $execute \ C_1$
JUMP h
 $g : execute \ C_2$
 $h :$

where g and h are **fresh** names.

G53CMP: Lecture 12 & 13 – p.24/46

Exercise: Code Function *execute*

Given

```
evaluate [[ I ]] =      addr(a) = [SB + 11]
  LOAD addr(I)        addr(b) = [SB + 12]
execute [[ I := IL ]] = addr(c) = [SB + 13]
  LOADL IL
  STORE addr(I)
```

generate code for:

```
if b then
  if c then a := 1 else a := 2
else
  a := 3
```

G53CMP: Lecture 12 & 13 – p.25/46

Code Function *execute* (5)

In detail (pseudo Haskell, code generation monad):

```
execute l env n [[ if E then C1 else C2 ]] = do
  g ← newName
  h ← newName
  evaluate l env E
  emit (JUMPIFZ g)
  execute l env n C1
  emit (JUMP h)
  emit (Label g)
  execute l env n C2
  emit (Label h)
```

G53CMP: Lecture 12 & 13 – p.26/46

Code Function *execute* (6)

```
execute [[ while E do C ]] =
  JUMP h
  g : execute C
  h : evaluate E
  JUMPIFNZ g
```

where *g* and *h* are **fresh** names.

G53CMP: Lecture 12 & 13 – p.27/46

Code Function *execute* (7)

```
execute [[ let Ds in C ]] =
  elaborate* Ds
  execute C
  POP 0 s
```

where *s* is the amount of storage allocated by *elaborate** *Ds*.

G53CMP: Lecture 12 & 13 – p.28/46

Code Function *execute* (8)

In detail (pseudo Haskell, code generation monad):

$$\begin{aligned} \text{execute } l \text{ env } n \llbracket \text{let } Ds \text{ in } C \rrbracket &= \mathbf{do} \\ & \quad (env', n') \leftarrow \text{elaborate}^* l \text{ env } n Ds \\ & \quad \text{execute } l \text{ env}' n' C \\ & \quad \text{emit } (\text{POP } 0 (n' - n)) \end{aligned}$$

where:

$$\begin{aligned} \text{elaborate}^* &: \text{Level} \rightarrow \text{CGEnv} \rightarrow \text{Depth} \\ & \quad \rightarrow \text{Declaration}^* \\ & \quad \rightarrow \text{CG TAMInst } (\text{Env}, \text{Depth}) \end{aligned}$$

G53CMP: Lecture 12 & 13 – p.29/46

Code Function *execute**

The code function *execute** has the obvious definition:

$$\text{execute}^* \llbracket \epsilon \rrbracket = \epsilon$$

$$\begin{aligned} \text{execute}^* \llbracket C ; Cs \rrbracket &= \\ & \quad \text{execute } C \\ & \quad \text{execute}^* Cs \end{aligned}$$

G53CMP: Lecture 12 & 13 – p.30/46

MiniTriangle Abstract Syntax Part II

<i>Expression</i>	→	<u>IntegerLiteral</u>	ExpLitInt
		<u>Identifier</u>	ExpVar
		<u>Operator</u> <i>Expression</i>	ExpUnOpApp
		<i>Expression</i> <u>Operator</u>	ExpBinOpApp
		<i>Expression</i>	
<i>Declaration</i>	→	const <u>Identifier</u> :	DeclConst
		<i>TypeDenoter</i> = <i>Expression</i>	
		var <u>Identifier</u> : <i>TypeDenoter</i>	DeclVar
		(:= <i>Expression</i> ϵ)	
<i>TypeDenoter</i>	→	<u>Identifier</u>	TDBaseType

G53CMP: Lecture 12 & 13 – p.31/46

Code Function *evaluate* (1)

$$\text{evaluate} : \text{Expression} \rightarrow \text{Instruction}^*$$

Fundamental invariant: all operations take arguments from the stack and writes result back onto the stack.

G53CMP: Lecture 12 & 13 – p.32/46

Code Function *evaluate* (2)

Consider evaluating $2 + 4 * 3 - 5$. Plausible instruction sequence:

LOADL 2	Stack: 2
LOADL 4	Stack: 4, 2
LOADL 3	Stack: 3, 4, 2
CALL mul	Stack: 12, 2
CALL add	Stack: 14
LOADL 5	Stack: 5, 14
CALL sub	Stack: 9

(mul, add, sub are routines in the MiniTriangle standard library.)

G53CMP: Lecture 12 & 13 – p.33/46

Code Function *evaluate* (3)

$$\text{evaluate } \llbracket IL \rrbracket =$$

LOADL c

where c is the value of IL .

$$\text{evaluate } \llbracket I \rrbracket =$$

fetch I

G53CMP: Lecture 12 & 13 – p.34/46

Code Function *evaluate* (4)

$$\text{evaluate } \llbracket \ominus E \rrbracket =$$

evaluate E
CALL *addr*(\ominus)

$$\text{evaluate } \llbracket E_1 \otimes E_2 \rrbracket =$$

evaluate E_1
evaluate E_2
CALL *addr*(\otimes)

(A call to a known function that can be replaced by a short code sequence can be optimised away at a later stage; e.g. CALL add \Rightarrow ADD.)

G53CMP: Lecture 12 & 13 – p.35/46

Code Functions *fetch* and *assign* (1)

In simplified MiniTriangle, all constants and variables are **global**. Hence addressing relative to SB.

$$\text{fetch } \llbracket I \rrbracket =$$

LOAD [SB + d]

where d is offset (or *displacement*) of I relative to SB.

$$\text{assign } \llbracket I \rrbracket =$$

STORE [SB + d]

where d is offset of I relative to SB.

G53CMP: Lecture 12 & 13 – p.36/46

Code Functions *fetch* and *assign* (2)

In a more realistic language, *fetch* and *assign* would take the current scope level and the scope level of the variable into account:

- Global variables addressed relative to SB.
- Local variables addressed relative to LB.
- Non-global variables in enclosing scopes would be reached by following the **static links** (see later lecture) in one or more steps, and *fetch* and *assign* would have to generate the appropriate code.

G53CMP: Lecture 12 & 13 – p.37/46

Exercise: Code Function *evaluate*

Given

$$\text{addr}(a) = [\text{SB} + 11]$$

$$\text{addr}(b) = [\text{SB} + 12]$$

$$\text{addr}(+) = \text{add}$$

$$\text{addr}(\ast) = \text{mult}$$

generate code for:

$$a + (b \ast 2)$$

G53CMP: Lecture 12 & 13 – p.39/46

Assignment revisited

In detail (pseudo Haskell, code generation monad) the code for assignment looks more like this.

Note that the variable actually is represented by an **expression** that gets evaluated to an **address**:

```
execute l env n [[ Ev := E ]] = do
  evaluate l env E
  evaluate l env Ev
  case sizeof(E) of
    1 → emit (STOREI 0)
    s → emit (STOREIB s)
```

(Reasons include: array references ($a[i]$), call by reference parameters.)

G53CMP: Lecture 12 & 13 – p.38/46

Code Function *elaborate* (1)

Elaboration must deposit value/reserve space for value on stack. Also, address (offset) of elaborated entity must be recorded (to be used by *fetch* and *assign*).

```
elaborate : Declaration → Instruction*
elaborate [[ const I : TD = E ]] =
  evaluate E
```

(Additionally, the offset (w.r.t. SB) has to be recorded for the identifier denoted by I .)

G53CMP: Lecture 12 & 13 – p.40/46

Code Function *elaborate* (2)

$$\begin{aligned} \text{elaborate } \llbracket \text{var } I : TD \rrbracket &= \\ &\quad \text{LOADL } 0 \\ \text{elaborate } \llbracket \text{var } I : TD := E \rrbracket &= \\ &\quad \text{evaluate } E \end{aligned}$$

(Additionally, the offset (w.r.t. SB) has to be recorded for the identifier denoted by I .)

LOADL 0 is just used to reserve space on the stack; the value of the literal does not matter. More space must be reserved if the values of the type are big (e.g. record, array).

G53CMP: Lecture 12 & 13 – p.41/46

Identifiers vs. Symbols (1)

- The coursework compiler HMTc uses **symbols** instead of identifiers in the latter stages.
- Symbols are introduced in the type checker (responsible for identification in HMTc) in place of identifiers (rep. changed from AST to MTIR).
- Symbols carry **semantic information** (e.g., type, scope level) to make that information readily available to e.g. the code generator.

(Cf. the lectures on identification where applied identifier occurrences were annotated with semantic information.)

G53CMP: Lecture 12 & 13 – p.43/46

Code Function *elaborate* (3)

For procedures and functions:

- Generate a fresh name for the entry point.
- Extend the environment according to formal argument declarations (the caller will push actual arguments onto stack prior to call).
- Generate code for the body at a scope level incremented by 1 and in the extended environment.

G53CMP: Lecture 12 & 13 – p.42/46

Identifiers vs. Symbols (2)

- Two kinds of (term-level) symbols:
 - External: defined outside the current compilation unit (e.g., in a library).
 - Internal: defined in the current compilation unit (in a `let`).

```
type TermSym = Either ExtTermSym IntTermSym
```

```
data ExtTermSym = ExtTermSym { ... }
```

```
data IntTermSym = IntTermSym { ... }
```

G53CMP: Lecture 12 & 13 – p.44/46

External Symbols

- External symbols are known entities.
- Can thus be looked up once and for all (during identification).
- Have a value, such as a (symbolic) address.

```
data ExtTermSym = ExtTermSym {
    etmsName :: Name,
    etmsType :: Type,
    etmsVal   :: ExtSymVal
}
```

```
data ExtSymVal = ESVLbl Name | ESVInt MTInt | ...
```

Internal Symbols

- Internal symbols do **not** carry any value such as stack displacement because this is not computed until the time of code generation.
- Such “late” information about an entity referred to by an internal symbol thus has to be looked up in the code generation environment.

```
data IntTermSym = IntTermSym {
    itmsLvl   :: ScopeLvl,
    itmsName  :: Name,
    itmsType  :: Type,
    itmsSrcPos :: SrcPos
}
```