

# Code Generation

## Compiler Design

### CSE 504

- 1 Syntax-Directed Code Generation
- 2 Machines
- 3 Expressions
- 4 Statements
- 5 Short-Circuit Code

Last modified: Wed Apr 08 2015 at 16:10:34 EDT  
Version: 1.6 15:28:43 2015/01/25  
Compiled at 16:12 on 2015/04/08



# Code Generation

- Intermediate code generation: Abstract (machine independent) code.
- Code optimization: Transformations to the code to improve time/space performance.
- Final code generation: Emitting machine instructions.

# Syntax Directed Translation

## Interpretation:

$$E \longrightarrow E_1 + E_2 \quad \{ E.val := E_1.val + E_2.val; \}$$

## Type Checking:

$$E \longrightarrow E_1 + E_2 \quad \left\{ \begin{array}{l} \text{if } E_1.type \equiv E_2.type \equiv int \\ \quad E.type = int; \\ \text{else} \\ \quad E.type = float; \\ \end{array} \right\}$$

# Code Generation via Syntax Directed Translation

## Code Generation:

$$E \rightarrow E_1 + E_2 \quad \left\{ \begin{array}{l} E.code = E_1.code \parallel \\ \quad E_2.code \parallel \\ \quad \text{"add"} \end{array} \right.$$

# Stack Machines

Simplified translation, but fewer opportunities for optimization.

- Machine Configuration:
  - Contents of stack; each element of the stack is a cell of some standard size (e.g. 32 bits).
  - Registers: Program counter, Stack pointer, (more later)
- Stack representation:
  - “[ ]” to represent empty stack
  - “ $v_1::S$ ” to represent a stack whose top element is  $v_1$  and the remainder of the stack is  $S$ .
  - “ $S[i]$ ” represents the value at the  $i$ -th element of stack  $S$  (counting from the base, not top, of the stack).

# Stack Machine Instructions

- `load_immed v`: Push  $v$  on stack.

$$S \quad \text{---} \boxed{\text{load\_immed } v} \quad \text{---} \quad v :: S$$

- `load`: Load value from given address to top of stack.

$$a :: S \quad \text{---} \boxed{\text{load}} \quad \text{---} \quad v :: S$$

where  $v = S[a]$ .

- `store`: Store a given value to a given address.

$$v :: a :: S \quad \text{---} \boxed{\text{store}} \quad \text{---} \quad v :: T$$

where  $T$  is same as  $S$  except  $T[a] = v$ .

- `add`: Add top two elements.

$$v_2 :: v_1 :: S \quad \text{---} \boxed{\text{add}} \quad \text{---} \quad v_1 + v_2 :: S$$

- `pop`: Remove top-most element.

$$v :: S \quad \text{---} \boxed{\text{pop}} \quad \text{---} \quad S$$

# Register Machines

- Machines with a (possibly unbounded) number of registers.  
Lets call them  $t_1, t_2, \dots$
- Separate *Heap* space for dynamically allocated objects.
- Instructions:
  - `move  $t_1, t_2$` : Move value from register  $t_2$  to  $t_1$ .
  - `move_immed  $t_1, i$` : move literal constant  $i$  to a register  $t_1$ .
  - `add  $t_1, t_2$` : Add values of  $t_1$  and  $t_2$ , store it back in  $t_1$ .

# Code Generation and Attributes (Stack Machine)

$$\begin{aligned}
 E &\longrightarrow E_1 + E_2 && \{ \\
 &&& E.code = E_1.code \parallel E_2.code \parallel \\
 &&& \quad \text{"add"} \\
 &&& \} \\
 E &\longrightarrow \text{int} && \{ E.code = \text{"load\_immed int.val"} \} \\
 E &\longrightarrow \text{id} && \{ E.code = \text{"load id.addr"} \} \\
 E &\longrightarrow \text{id} = E_1 && \{ \\
 &&& E.code = \text{"load\_immed id.addr"} \parallel \\
 &&& \quad E_1.code \parallel \\
 &&& \quad \text{"store"} \\
 &&& \}
 \end{aligned}$$


---

`id.addr` is the address of cell allocated for `id`.



# Variables and Addresses in Register Machines

- Since the register machine has unbounded number of registers, all local variables and function parameters will be stored in registers.
- Each program variable is mapped to a distinct abstract register. Each `id` has an attribute `id.addr` to represent this mapping.
- For final code generation, the registers of the abstract machine will be mapped to (a small, finite) set of registers of a concrete machine (e.g. MIPS).
- Clearly not all abstract machine registers may have a corresponding concrete machine register. Such abstract registers will be **spilled** to cells on a **stack** on the concrete machine.

# Code Generation and Attributes (Register Machine)

```
 $E \rightarrow E_1 + E_2$  {  
     $E.t = \text{generate\_new\_temporary}();$   
     $E.code = E_1.code \parallel E_2.code$   
         $\parallel \text{"add } E.t, E_1.t, E_2.t"$   
}  
 $E \rightarrow \text{int}$  {  
     $E.t = \text{generate\_new\_temporary}();$   
     $E.code = [\text{"mov\_immed } E.t, \text{int.val"}]$   
}  
 $E \rightarrow \text{id}$  {  
     $E.t = \text{id.addr};$   
     $E.code = []$   
}  
 $E \rightarrow \text{id} = E_1$  {  
     $E.t = E_1.t;$   
     $E.code = E_1.code$   
         $\parallel \text{"move id.addr, } E.t"$   
}
```

## References

Some languages allow variables to have *locations* (i.e. addresses) and a variable may **refer** to another variable's location.

- Languages such as C/C++ permit programmers to obtain location of arbitrary variables (using the “address-of” operation, “&”), and dereference locations (i.e. access the value stored at an address, using “\*”).
- When translating a C-like language, every variable should be potentially allocated on stack; it can be mapped to a register if there is no operation that takes its address.
- Languages such as Java give locations only to objects and arrays. All variables are stack-allocated.  
The location of a variable itself is not accessible to the program.

## *l*- and *r*-values

```
i = i + 1;
```

- ***r*-value**: actual value of the expression
- ***l*-value**: for expressions associated with specific memory addresses, the location where the value of the expression is stored.
- Some expressions (e.g. `i`) have both *l*- and *r*-values.
- Other expressions (e.g. `5`) have only an *r*-value.
- Some expression's values may be (at least temporarily) stored in locations, but those locations may not have a meaning in terms of the program, and we consider them also to have only *r*-values (e.g. `i+1`).

## Code Generation for $L$ -expressions

An  $L$ -expression is one which has an  $l$ -value.

- Roughly speaking  $L$ -expressions are those that may occur on the lhs of an assignment.
- In Proto(2), the only  $L$ -expressions were identifiers.
- In Proto(3),  $L$ -expressions include array access.
- For compiling assignments, we will use additional attributes for  $L$ -expressions (other than  $L.t$  and  $L.code$ , which all expressions have).

# Arrays

Consider expression grammar changed as follows:

$$E \rightarrow E + E$$
$$E \rightarrow L = E$$
$$E \rightarrow \text{int}$$
$$E \rightarrow L$$
$$L \rightarrow \text{id}$$
$$L \rightarrow L[E]$$

- $L$  represents simple identifiers as well as array expressions.
- The index of an array expression can be any arbitrary expression (including an array expression itself)  
Example:  $x[y[i]]$
- The base of an array expression is an identifier or another array expression.  
Example:  $(x[i])[j]$
- LHS of an assignment can be an array expression.

# Addresses and Allocation

- For Proto, we'll use Java-like convention of keeping variables in stack/registers, and arrays (and later, objects) on heap.
- For heap access, we use the following intermediate code instructions:
  - `hstore a, r`: store value to a heap cell.  
Register  $a$  has the address of the cell in heap, and register  $r$  has the value to be stored.
  - `hload r, a`: load value from a heap cell.  
Register  $a$  has the address of the cell in heap, and register  $r$  is the destination for the load.
  - `halloc r1, r2`: allocate a segment of heap cells.  
Register  $r_2$  contains the number of cells to allocate. Register  $r_1$  will then be set to the base address of the allocated heap cells.
  - ?? `hsize r1, r2`: get size of a heap segment. Bounds Check  
Register  $r_1$  is the address of the heap segment. Register  $r_2$  will then be set to the size of the segment.

## Generating code for arrays:

```
 $E \rightarrow \text{new } T [ E_1 ]$   
{  
   $E.t = \text{generate\_new\_temporary}();$   
   $E.code = E_1.code$   
    || "halloc  $E.t, E_1.t$ "  
}
```

## Allocation

- $E_1$  will be an integer-valued expression that specifies the number of elements in the array to allocate.
- Type  $T$  is ignored (at least, for now).
- $E$ , then, is a reference to the newly allocated array.
- If bounds check is needed, additional book-keeping info needs to be maintained with the array.  
... allocate  $n + 1$  cells, and use the zero-th cell to store the length!



## Generating code for arrays:

## LHS

```

L  →  id      {
                L.t = L.at = id.addr;
                L.lcode = L.rcode = [ ];
                L.mem = reg;
            }

L  →  L1 [ E ] {
                L.at = generate_new_temporary();
                L.lcode = L1.rcode
                    || E.code
                    || "mul L.at, E.t, 4"
                    || "add L.at, L.at, L1.t";
                L.t = generate_new_temporary();
                L.rcode = L.lcode
                    || "hload L.t, L.at";
                L.mem = heap;
            }

```

- *L.t*: Register holding *L*'s value.
- *L.at*: Register holding *L*'s address.
- *L.lcode*: Code for evaluating *L*'s address.
- *L.rcode*: Code for evaluating *L*'s value.
- Note: no bounds check!

## Generating code for arrays:

## RHS

$$E \longrightarrow L \left\{ \begin{array}{l} E.t = L.t; \\ E.code = L.rcode \end{array} \right\}$$

Example expression:

```
(i + a[i])  
  + b[i][j]
```

With  $i.addr = t_1$ ,  
 $j.addr = t_2$ ,  
 $a.addr = t_3$ ,  
 $b.addr = t_4$ .

```
// i's rcode (empty)  
// a[i]'s rcode  
mul t5, t1, 4  
add t5, t5, t3  
aload t6, t5  
// i+a[i]'s code  
add t7, t1, t6  
// b[i][j]'s rcode:  
//   b[i]'s rcode  
mul t8, t1, 4  
add t8, t8, t4  
aload t9, t8  
//   use b[i] as base:  
mul t10, t2, 4  
add t10, t10, t8  
aload t11, t10  
// add b[i][j] to prev result  
add t12, t7, t11
```

## Generating code for arrays:

## Assignments

```

E  →  L = E1
      {
        E.t = E1.t;
        if L.mem == reg
            assigncode = "move L.at, E1.t";
        else
            assigncode = "hstore L.at, E1.t";
        E.code = L.lcode
            || E1.code
            || assigncode;
      }

```

# Code Generation for Statements

$$\begin{aligned} Ss &\longrightarrow S Ss_1 && \left\{ \begin{array}{l} Ss.code = S.code \\ \quad \parallel Ss_1.code; \end{array} \right. \\ & && \left. \right\} \\ Ss &\longrightarrow \epsilon && \left\{ Ss.code = [ ] \right\} \\ S &\longrightarrow E ; && \left\{ \begin{array}{l} S.code = E.code; \end{array} \right. \\ & && \left. \right\} \end{aligned}$$

# Conditional Statements

```
S  →  if E, S1, S2  {  
      elselabel = get_new_label();  
      endlabel = get_new_label();  
      S.code =      E.code  
                  ||      "beq E.t, 0, elselabel"  
                  ||      S1.code;  
                  ||      "jmp endlabel"  
                  || "elselabel:"  
                  ||      S2.code;  
                  || "endlabel:"  
      }
```

# Conditional Statements and Continuations

$S.end$ : label to jump after  $S$  is executed completely.

```
 $S \rightarrow$  if  $E, S_1, S_2$  {  
     $S.begin = get\_new\_label();$   
     $S_1.end = S_2.end = S.end;$   
     $S.code = "S.begin:"$   
         $\| E.code$   
         $\| "beq E.t, 0, S_2.begin"$   
         $\| S_1.code$   
         $\| S_2.code;$   
}
```

# Continuations

Attributes of a statement that specify where control will flow to after the statement is executed.

- Analogous to the *follow* sets of grammar symbols.
- In deterministic languages, there is only one continuation for each statement.
- Can be generalized to include local variables whose values are needed to execute the following statements:

*Uniformly captures call, return and exceptions.*

## Sequence and Continuation

- Most frequently, the continuation of a statement will simply be its succeeding statement.
- We will use a special label “*fallthrough*” to denote this.

$$Ss \longrightarrow S Ss_1 \quad \left\{ \begin{array}{l} Ss_1.end = Ss.end; \\ S.end = \textit{fallthrough}; \\ Ss.code = \dots \end{array} \right\}$$
$$S \longrightarrow E ; \quad \left\{ \begin{array}{l} \textit{if } S.end == \textit{fallthrough} \\ \quad \textit{next} = [ ] \\ \textit{else} \\ \quad \textit{next} = \textit{"jmp } S.end\textit{"} \\ S.code = \dots \\ \quad \parallel \textit{next}; \end{array} \right\}$$



## Code Generation for Boolean Expressions

$$E \longrightarrow E_1 \ \&\& \ E_2 \quad \left\{ \begin{array}{l} E.t = \text{generate\_new\_temporary}(); \\ E.code = E_1.code \\ \quad \parallel E_2.code \\ \quad \parallel \text{"and } E.t, E_1.t, E_2.t\text{"}; \end{array} \right\}$$

- The above code evaluates  $E_2$  regardless of the value of  $E_1$ .
- Short circuit code: evaluate  $E_2$  only if needed.

$$E \longrightarrow E_1 \ \&\& \ E_2 \quad \left\{ \begin{array}{l} E.t = \text{generate\_new\_temporary}(); \\ \text{skip} = \text{generate\_new\_label}(); \\ E.code = E_1.code \\ \quad \parallel \text{"move } E.t, E_1.t\text{"} \\ \quad \parallel \text{"beq } E.t, 0, \text{skip}\text{"} \\ \quad \parallel E_2.code \\ \quad \parallel \text{"move } E.t, E_2.t\text{"} \\ \quad \parallel \text{"skip:"}; \end{array} \right\}$$

## Generating Shortcircuit Code

Use two continuations for each boolean expression:

- *E.success*: where control will go when expression in *E* evaluates to *true*.
- *E.fail*: where control will go when expression in *E* evaluates to *false*.

Both continuations are *inherited* attributes.

# Shortcircuit Code for Boolean Expressions

$$\begin{aligned} E &\longrightarrow E_1 \ \&\& \ E_2 \quad \left\{ \begin{array}{l} E_1.fail = E.fail; \\ E_2.fail = E.fail; \\ E_1.success = get\_new\_label(); \\ E_2.success = E.success; \\ E.code = E_1.code \parallel \\ \qquad \qquad \qquad \text{"E}_1.success\text{"} \parallel \\ \qquad \qquad \qquad E_2.code \end{array} \right\} \\ E &\longrightarrow ! E_1 \quad \left\{ \begin{array}{l} E_1.fail = E.success; \\ E_1.success = E.fail; \\ E.code = E_1.code \end{array} \right\} \\ E &\longrightarrow \text{true} \quad \left\{ \begin{array}{l} E.code = \text{"jmp, E.success"} \end{array} \right\} \end{aligned}$$

# Short-circuit code for Conditional Statements

```

S  →  if E, S1, S2  {
      S.begin = get_new_label();
      S1.end = S2.end = S.end;
      E.success = S1.begin;
      E.fail = S2.begin;
      S.code = "S.begin:" ||
               E.code ||
               S1.code ||
               S2.code;
    }
  
```

# Continuations and Code Generation

Continuation of a statement is an inherited attribute.

*It is not an L-inherited attribute!*

Code of statement is a synthesized attribute, but is dependent on its continuation.

**Backpatching:** *Make two passes to generate code.*

- 1 Generate code, leaving “holes” where continuation values are needed.
- 2 Fill these holes on the next pass.

# What's left?

After intermediate code is generated,

- **Optimize** intermediate code using target machine-independent techniques.

Examples:

- constant propagation
  - loop-invariant code motion
  - dead-code elimination
  - strength reduction
- **Generate** final machine code  
Perform target machine-specific optimizations.