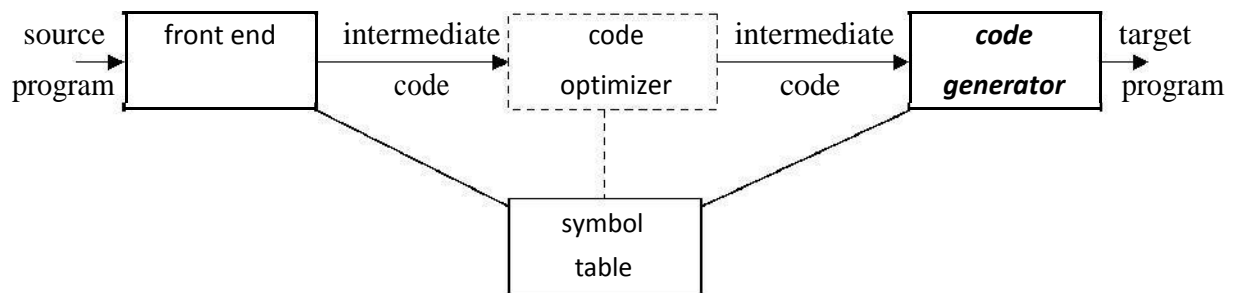


UNIT V

CODE GENERATION

The final phase in compiler model is the code generator. It takes as input an intermediate representation of the source program and produces as output an equivalent target program. The code generation techniques presented below can be used whether or not an optimizing phase occurs before code generation.

Position of code generator



ISSUES IN THE DESIGN OF A CODE GENERATOR

The following issues arise during the code generation phase :

- Input to code generator
- Target program
- Memory management
- Instruction selection
- Register allocation
- Evaluation order

Input to code generator:

The input to the code generation consists of the intermediate representation of the source program produced by front end , together with information in the symbol table to determine run-time addresses of the data objects denoted by the names in the intermediate representation.

Intermediate representation can be :

- Linear representation such as postfix notation
- Three address representation such as quadruples
- Virtual machine representation such as stack machine code
- Graphical representations such as syntax trees and dags.

Prior to code generation, the front end must be scanned, parsed and translated into intermediate representation along with necessary type checking. Therefore, input to code generation is assumed to be error-free.

Target program:

The output of the code generator is the target program. The output may be :

- a. Absolute machine language

It can be placed in a fixed memory location and can be executed immediately.

Relocatable machine language

It allows subprograms to be compiled separately.

Assembly language

Code generation is made easier.

Memory management:

Names in the source program are mapped to addresses of data objects in run-time memory by the front end and code generator.

It makes use of symbol table, that is, a name in a three-address statement refers to a symbol-table entry for the name.

Labels in three-address statements have to be converted to addresses of instructions. For example,

j : **goto** i generates jump instruction as follows :

- if $i < j$, a backward jump instruction with target address equal to location of code for quadruple i is generated.
- if $i > j$, the jump is forward. We must store on a list for quadruple i the location of the first machine instruction generated for quadruple j . When i is processed, the machine locations for all instructions that forward jumps to i are filled.

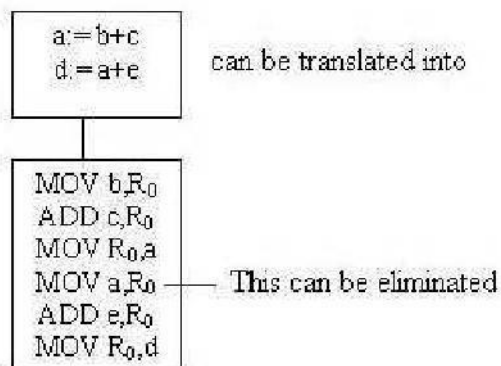
Instruction selection:

The instructions of target machine should be complete and uniform.

Instruction speeds and machine idioms are important factors when efficiency of target program is considered.

- The quality of the generated code is determined by its speed and size.

The former statement can be translated into the latter statement as shown below:



Register allocation

Instructions involving register operands are shorter and faster than those involving operands in memory.

The use of registers is subdivided into two subproblems :

Register allocation – the set of variables that will reside in registers at a point in the program is selected.

- **Register assignment** – the specific register that a variable will reside in is picked.

Certain machine requires even-odd *register pairs* for some operands and results. For example, consider the division instruction of the form :

D x, y

where, x – dividend even register in even/odd register

pair y – divisor

even register holds the remainder

odd register holds the quotient

Evaluation order

The order in which the computations are performed can affect the efficiency of the target code. Some computation orders require fewer registers to hold intermediate results than others.

TARGET MACHINE

Familiarity with the target machine and its instruction set is a prerequisite for designing a good code generator.

The target computer is a byte-addressable machine with 4 bytes to a word.

It has n general-purpose registers, R_0, R_1, \dots, R_{n-1} .

It has two-address instructions of the form:

op source, destination

where, *op* is an op-code, and *source* and *destination* are data fields.

It has the following op-codes :

MOV (move *source* to *destination*)

ADD (add *source* to *destination*)

SUB (subtract *source* from *destination*)

The *source* and *destination* of an instruction are specified by combining registers and memory locations with address modes.

Address modes with their assembly-language forms

MODE	FORM	ADDRESS	ADDED COST
<i>Absolute</i>	M	M	1
<i>Register</i>	R	R	0
<i>indexed</i>	$c(R)$	$c + \text{contents}(R)$	1
<i>indirect register</i>	*R	$\text{contents}(R)$	0
<i>indirect indexed</i>	* $c(R)$	$\text{contents}(c + \text{contents}(R))$	1
<i>literal</i>	# c	c	1

For example : MOV R0, M stores contents of Register R0 into memory location M ;
MOV 4(R0), M stores the value *contents(4+contents (R0))* into M.

Instruction costs :

Instruction cost = 1+cost for source and destination address modes. This cost corresponds to the length of the instruction.

Address modes involving registers have cost zero.

Address modes involving memory location or literal have cost one.

Instruction length should be minimized if space is important. Doing so also minimizes the time taken to fetch and perform the instruction.

For example : MOV R0, R1 copies the contents of register R0 into R1. It has cost one, since it occupies only one word of memory.

The three-address statement **a := b + c** can be implemented by many different instruction sequences :

i) MOV b, R0

ADD c, R0 cost = 6

MOV R0, a

ii) MOV b, a

ADD c, a cost = 6

iii) Assuming R0, R1 and R2 contain the addresses of a, b, and c :

MOV *R1, *R0

ADD *R2, *R0 cost = 2

In order to generate good code for target machine, we must utilize its addressing capabilities efficiently.

RUN-TIME STORAGE MANAGEMENT

Information needed during an execution of a procedure is kept in a block of storage called an activation record, which includes storage for names local to the procedure.

The two standard storage allocation strategies are:

Static allocation

Stack allocation

In static allocation, the position of an activation record in memory is fixed at compile time.

In stack allocation, a new activation record is pushed onto the stack for each execution of a procedure. The record is popped when the activation ends.

The following three-address statements are associated with the run-time allocation and deallocation of activation records:

Call,

Return,

Halt, and

Action, a placeholder for other statements.

We assume that the run-time memory is divided into areas for:

Code

Static data

Stack

Static allocation

Implementation of call statement:

The codes needed to implement static allocation are as follows:

```
MOV #here +20, callee.static_area      /*It saves return address*/
```

```
GOTO callee.code_area      /*It transfers control to the target code for the called procedure */
```

where,

callee.static_area – Address of the activation record *callee.code_area*

– Address of the first instruction for called procedure

#here +20 – Literal return address which is the address of the instruction following GOTO.

Implementation of return statement:

A return from procedure *callee* is implemented by :

```
GOTO * callee.static_area
```

This transfers control to the address saved at the beginning of the activation record.

Implementation of action statement:

The instruction ACTION is used to implement action statement.

Implementation of halt statement:

The statement HALT is the final instruction that returns control to the operating system.

Stack allocation

Static allocation can become stack allocation by using relative addresses for storage in activation records. In stack allocation, the position of activation record is stored in register so words in activation records can be accessed as offsets from the value in this register.

The codes needed to implement stack allocation are as follows:

Initialization of stack:

```
MOV #stackstart , SP      /* initializes stack */
```

Code for the first procedure

```
HALT                      /* terminate execution */
```

Implementation of Call statement:

```
ADD #caller.recordsize, SP  /* increment stack pointer */
```

```
MOV #here +16, *SP        /*Save return address */
```

```
GOTO callee.code_area
```

where,
caller.recordsize – size of the activation record
#here +16 – address of the instruction following the **GOTO**

Implementation of Return statement:

GOTO *0 (SP) /*return to the caller */

SUB *#caller.recordsize*, SP /* decrement SP and restore to previous value */

A SIMPLE CODE GENERATOR

A code generator generates target code for a sequence of three- address statements and effectively uses registers to store operands of the statements.

For example: consider the three-address statement **a := b+c** It can have the following sequence of codes:

ADD R_j, R_i Cost = 1 // if R_i contains b and R_j contains c

(or)

ADD c, R_i Cost = 2 // if c is in a memory location

(or)

MOV c, R_j Cost = 3 // move c from memory to R_j and add

ADD R_j, R_i

Register and Address Descriptors:

A register descriptor is used to keep track of what is currently in each registers. The register descriptors show that initially all the registers are empty.

An address descriptor stores the location where the current value of the name can be found at run time.

A code-generation algorithm:

The algorithm takes as input a sequence of three -address statements constituting a basic block. For each three-address statement of the form **x := y op z**, perform the following actions:

Invoke a function *getreg* to determine the location L where the result of the computation y op z should be stored.

Consult the address descriptor for y to determine y', the current location of y. Prefer the register for y' if the value of y is currently both in memory and a register. If the value of y is not already in L, generate the instruction **MOV y' , L** to place a copy of y in L.

Generate the instruction **OP z' , L** where z' is a current location of z. Prefer a register to a memory location if z is in both. Update the address descriptor of x to indicate that x is in location L. If x is in L, update its descriptor and remove x from all other descriptors.

If the current values of y or z have no next uses, are not live on exit from the block, and are in registers, alter the register descriptor to indicate that, after execution of $x := y \text{ op } z$, those registers will no longer contain y or z.

Generating Code for Assignment Statements:

The assignment $d := (a-b) + (a-c) + (a-c)$ might be translated into the following three-address code sequence:

$t := a - b$

$u := a - c$

$v := t + u$

$d := v + u$

with d live at the end.

Code sequence for the example is:

Statements	Code Generated	Register descriptor	Address descriptor
		Register empty	
$t := a - b$	MOV a, R0 SUB b, R0	R0 contains t	t in R0
$u := a - c$	MOV a, R1 SUB c, R1	R0 contains t R1 contains u	t in R0 u in R1
$v := t + u$	ADD R1, R0	R0 contains v R1 contains u	u in R1 v in R0
$d := v + u$	ADD R1, R0 MOV R0, d	R0 contains d	d in R0 d in R0 and memory

Generating Code for Indexed Assignments

The table shows the code sequences generated for the indexed assignment statements **a := b [i]** and **a [i] := b**

Statements	Code Generated	Cost
a := b[i]	MOV b(Ri), R	2
a[i] := b	MOV b, a(Ri)	3

Generating Code for Pointer Assignments

The table shows the code sequences generated for the pointer assignments **a := *p** and ***p := a**

Statements	Code Generated	Cost
a := *p	MOV *Rp, a	2
*p := a	MOV a, *Rp	2

Generating Code for Conditional Statements

Statement	Code
if x < y goto z	CMP x, y CJ< z /* jump to z if condition code is negative */
x := y + z if x < 0 goto z	MOV y, R0 ADD z, R0 MOV R0, x CJ< z