

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:

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

Coercion

Type Checking:

$E \longrightarrow E_1 + E_2$ {
if $E_1.type \equiv E_2.type \equiv int$
 $E.type = int;$
else
 $E.type = float;$
}

Code Generation via Syntax Directed Translation

Print 3+(2*5)

Code Generation:

$E \rightarrow E_1 + E_2$ {
 $E.code = E_1.code ||$
 $E_2.code ||$
“add”)
}

Postfix code

Intermediate Code

“Abstract” code generated from AST

- **Simplicity and Portability**

- Machine independent code.
- Enables common optimizations on intermediate code.
- Machine-dependent code optimizations postponed to last phase.

Intermediate Forms

- Stack machine code:

Code for a “postfix” stack machine.

- Two address code:

Code of the form “add r_1, r_2 ”

$$r_2 = r_1 + r_2$$

- Three address code:

Code of the form “add $src_1, src_2, dest$ ”

$$dest = src_1 + src_2$$

Quadruples and Triples: Representations for three-address code.

Quadruples

Explicit representation of three-address code.

Example: $a := a + b * -c;$

Instr	Operation	Arg 1	Arg 2	Result
→ (0)	<u>uminus</u>	<u>c</u>		<u>t_1</u>
→ (1)	mult	b	t_1	t_2
→ (2)	add	a	t_2	t_3
(3)	move	t_3		a

Triples

Representation of three-address code with implicit destination argument.

Example: $a := a + b * -c;$

Instr	Operation	Arg 1	Arg 2
(0)	uminus	c	
(1)	mult	b	(0)
(2)	add	a	(1)
(3)	move	a	(2)

Intermediate Forms

Choice depends on convenience of further processing

- Stack code is simplest to generate for expressions.
- Quadruples are most general, permitting most optimizations including code motion.
- Triples permit optimizations such as *common subexpression elimination*, but code motion is difficult.

Static Single Assignment (SSA) ←

- Each variable is assigned at most once
- ϕ nodes used to combine values of variables after a conditional

```
if (f) x = 1; else x=2;  
y=x*x;
```

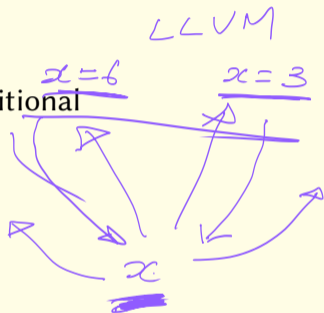
Becomes

```
if (f) x1 = 1; else x2=2;
```

```
x3 =  $\phi$ (x1, x2);
```

```
y=x3*x3;
```

merging



```
x1 = 3;  
if (c)
```

```
x = 5  
2
```

```
x3 =  $\phi$ (x1, x2)
```

Generating 3-address code

```
E → E1 + E2 {  
  E.addr = newtemp();  
  E.code = E1.code || E2.code ||  
           E.addr || ':' || E1.addr || '+' || E2.addr;  
}
```

```
E → int {  
  E.addr = newtemp();  
  E.code = E.addr || ':' || int.val;  
}
```

```
E → id {  
  E.addr = id.name;  
  E.code = ""; // no code needed  
}
```

```
⌈ 20  
 21  
 22  
⌋
```

E.code is a string-valued attribute containing the code for E. It generates a temp var for storing the result of evaluating E, &

→ string concatenation

generates code to evaluate E and store in this temp.

Generation of Postfix Code for Boolean Expressions

$E \rightarrow E_1 \ \&\& \ E_2 \{$
 $\quad E.code = E_1.code \ ||$
 $\quad \quad E_2.code \ ||$
 $\quad \quad gen(\text{and})$
 $\quad \}$

$E \rightarrow ! E_1 \{$
 $\quad E.code = E_1.code \ ||$
 $\quad \quad gen(\text{not})$
 $\quad \}$

$E \rightarrow \text{true} \{E.code = gen(\text{load_immed}, 1) \}$

$E \rightarrow \text{id} \{E.code = gen(\text{load}, \text{id.addr}) \}$

Code for Boolean Expressions

```
if ((p != NULL)
    && (p->next != q)) {
    ... then part
} else {
    ... else part
}
```

E's code →

```
load(p);
null();
neq();
load(p);
ildd(1);
getfield();
load(q);
neq();
and();
jnz elselabel;
... then part
elselabel:
... else part
```

E₁ code

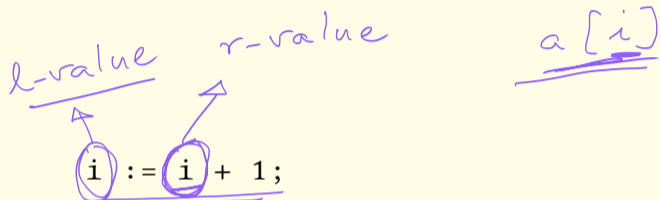
E₂ code

Shortcircuit Code

```
if ((p != NULL)
    && (p->next != q)) {
    ... then part
} else {
    ... else part
}
```

```
load(p);
null();
neq();
  jnz elselabel;
load(p);
ildc(1);
  getfield();
load(q);
neq();
  jnz elselabel;
  ... then part
elselabel:
  ... else part
```

l- and r-Values



- l-value: location where the value of the expression is stored.

- r-value: actual value of the expression (or

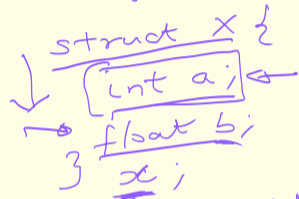
- for variables, the r-value is the value stored at the variable's location (which, in turn, is the l-value of the variable)

Computing *l*-values

$E \rightarrow \underline{id} \{$
 $\quad \underline{E.lval} = \underline{id.loc};$
 $\quad \underline{E.code} = ";"; \}$ // empty code

loc attribute stores location of the variable id.

$E \rightarrow \underline{E_1} [\underline{E_2}] \{$
 $\quad \underline{E.lval} = \underline{newtemp()};$
 $\quad \underline{x} = \underline{newtemp()};$
 $\quad \underline{E.lcode} = \underline{E_1.lcode} \parallel \underline{E_2.code} \parallel$
 $\quad \underline{x} \parallel \underline{':='} \parallel \underline{E_2.rval} \parallel \underline{ '*' } \parallel \underline{E_1.elemsize} \parallel$
 $\quad \underline{E.lval} \parallel \underline{':='} \parallel \underline{E_1.lval} \parallel \underline{ '+' } \parallel \underline{x} \}$



$x.a$ — resides at the same location as x .
 $x.b$ resides at location of $x + \text{size of (int)}$

$E \rightarrow \underline{E_1} . \underline{id} \{ // \text{for field access}$
 $\quad \underline{E.lval} = \underline{newtemp()};$
 $\quad \underline{E.lcode} = \underline{E_1.lcode} \parallel$
 $\quad \underline{E.lval} \parallel \underline{':='} \parallel \underline{E_1.lval} \parallel \underline{ '+' } \parallel \underline{id.offset} \}$

Computing lval and rval attributes

$$E \rightarrow \dot{E}_1 = \dot{E}_2 \{$$
$$E.code = E_1.lcode || E_2.code ||$$
$$\text{gen}(' * ' E_1.lval ':=' E_2.rval)$$
$$E.rval = E_2.rval \}$$
$$E \rightarrow E_1 [E_2] \{$$
$$E.lval = \text{newtemp}();$$
$$E.rval = \text{newtemp}();$$
$$x = \text{newtemp}();$$
$$E.lcode = E_1.lcode || E_2.code ||$$
$$\text{gen}(x ':=' E_2.rval * E_1.elemsize) ||$$
$$\text{gen}(E.lval ':=' E_1.lval '+' x)$$
$$E.code = E.lcode ||$$
$$\text{gen}(E.rval ':=' '*' E.lval)$$
$$\}$$

$*x = y$

store y, x

$e[d[i]] = \underline{5+3}$

Function Calls (Call-by-Value)

E \rightarrow $E_1(E_2, E_3)$ {
 $E.rval = newtemp();$
 $E.code = E_1.code ||$
 $E_2.code ||$
 $E_3.code ||$
 $gen(push E_2.rval)$
 $gen(push E_3.rval)$
 $gen(call E_1.rval)$
 $gen(pop E.rval)$
}

$(a \rightarrow f)(3, 5)$

also need to
pop arguments

Function Calls (Call-by-Reference)

$E \longrightarrow E_1 (\underline{E_2}, \underline{E_3}) \{$

$E.rval = newtemp();$

$E.code = E_1.code \parallel$

$E_2.lcode \parallel$

$E_3.lcode \parallel$

$gen(push E_2.lval)$

$gen(push E_3.lval)$

$gen(call E_1.rval)$

$gen(pop E.rval)$

$\}$

Code Generation for Statements

$$S \longrightarrow \underline{S_1 ; S_2} \left\{ \begin{array}{l} S.code = \underline{S_1.code} \parallel \\ \quad \underline{S_2.code}; \end{array} \right.$$
$$\underline{S} \longrightarrow E \quad \left\{ \underline{S.code = E.code}; \right.$$

Conditional Statements

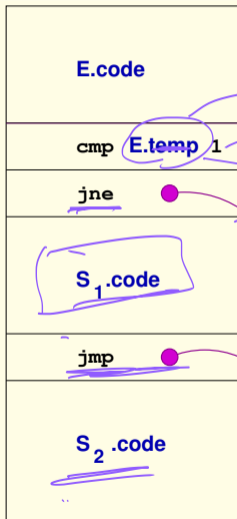
$S \rightarrow \text{if } E, S_1, S_2$

then

else

S.elsepart:

S.end:



represents "true"
previously called E.addr

Conditional Statements

$S \rightarrow$ if E , S_1 , S_2 {

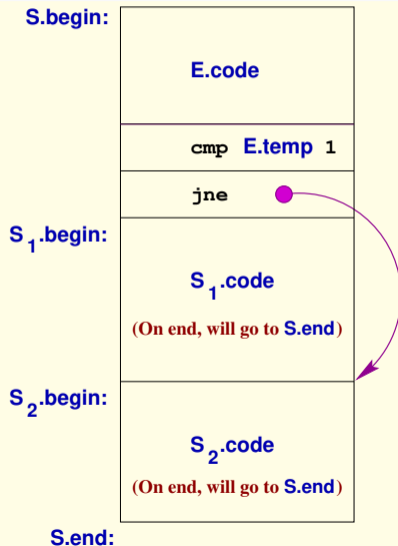
$elselabel = newlabel();$
 $endlabel = newlabel();$
 $S.code =$ $E.code ||$
 $gen(if\ E.temp\ \neq\ '1'\ elselabel) ||$
 $S_1.code ||$
 $gen(jmp\ endlabel) ||$
 $gen(elselabel:) ||$
 $S_2.code ||$
 $gen(endlabel:)$
}

L1
L2
L3
⋮

"if" ||
 $E.temp$ ||
"≠" ||
 $elselabel$

If Statements: An Alternative

$S \longrightarrow \text{if } E, S_1, S_2$



Continuations

An attribute of a statement that specifies 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**.

Conditional Statements and Continuations

```
S → if(E, S1, S2) {  
    S.begin = newlabel();  
    S.end = newlabel();  
    S1.end = S2.end = S.end;  
    S.code = gen(S.begin:) ||  
              E.code ||  
              gen(if E.rval '==' '1' S2.begin) ||  
              S1.code ||  
              S2.code ||  
              gen(S.end:)  
}
```

Continuations

- Each boolean expression has two possible continuations:
 - *E.true*: where control will go when expression in *E* evaluates to *true*.
 - *E.false*: where control will go when expression in *E* evaluates to *false*.
- Every statement *S* has one continuation, *S.next*
- Every while loop statement has an additional continuation, *S.begin*

Shortcircuit Code for Boolean Expressions

$E \rightarrow E_1 \ \&\& \ E_2 \{$
 $E_1.true = newlabel();$
 $E_1.false = E_2.false = E.false;$
 $E_2.true = E.true;$
 $E.code = E_1.code \ || \ gen(E_1.true ':') \ || \ E_2.code$
 $\}$

$E \rightarrow$ comparison

$E \rightarrow E_1 \ \text{or} \ E_2 \{$
 $E_1.true = E_2.true = E.true;$
 $E_1.false = newlabel();$
 $E_2.false = E.false;$
 $E.code = E_1.code \ || \ gen(E_1.false ':') \ || \ E_2.code$
 $\}$

$E \rightarrow ! E_1 \{$
 $E_1.false = E.true; E_1.true = E.false;$
 $\}$

$E \rightarrow true \{ E.code = gen(jmp, E.true) \}$

Short-circuit code for Conditional Statements

```
S → S1; S2 {  
  S1.next = newlabel();  
  S.code = S1.code || gen(S1.next ':') || S2.code;  
}
```

```
S → if E then S1 else S2 {  
  E.true = newlabel();  
  E.false = newlabel();  
  S1.next = S2.next = S.next;  
  S.code = E.code ||  
    gen(E.true ':') || S1.code ||  
    gen(jmp S.next) ||  
    gen(E.false ':') || S2.code;  
}
```

Continuation-style
code for if-then-else

Short-circuit code for While

```
S → while E do S1 {  
  S.begin = newlabel();  
  E.true = newlabel();  
  E.false = S.next;  
  S1.next = S.begin;  
  S.code = gen(S.begin':') || E.code ||  
           gen(E.true':') || S1.code ||  
           gen(jmp S.begin);  
}
```


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.

Machine Code Generation Issues

- Register assignment
- Instruction selection
- ...

How GCC Handles Machine Code Generation

- gcc uses machine descriptions to *automatically* generate code for target machine
- machine descriptions specify:
 - memory addressing (bit, byte, word, big-endian, ...)
 - registers (how many, whether general purpose or not, ...)
 - stack layout
 - parameter passing conventions
 - semantics of instructions
 - ...

Specifying Instruction Semantics

- gcc uses intermediate code called RTL, which uses a LISP-like syntax
- after parsing, programs are translated into RTL
- semantics of each instruction is also specified using RTL:

```
movl (r3), @8(r4) ≡
```

```
  (set (mem: SI (plus: SI (reg: SI 4) (const_int 8)))  
       (mem: SI (reg: SI 3)))
```

- cost of machine instructions also specified
- gcc code generation = selecting a low-cost instruction sequence that has the same semantics as the intermediate code