CS 4300: Compiler Theory

Chapter 6 Intermediate-Code Generation

Xuejun Liang 2024 Fall

Introduction

Logical structure of a compiler front end



Syntax trees are high level

Three-address code can range from high-level to low-level, depending on the choice of operators

Static versus Dynamic Checking

- Static checking: checked at compile time
 - Compiler enforces programming language's static semantics
 - Typical examples of static checking:
 - Type checks
 - Flow-of-control checks
 - Uniqueness checks
 - Name-related checks
- Dynamic semantics: checked at run time
 - Compiler generates verification code to enforce programming language's dynamic semantics

Type Checking, Overloading, Coercion, Polymorphism

```
class X { virtual int m(); } *x;
class Y: public X { virtual int m(); } *y;
int op(int), op(float);
int f(float);
int a, c[10], d;
d = c + d; // FAIL
*d = a;
             // FAIL
a = op(d); // OK: static overloading (C++)
a = f(d); // OK: coersion of d to float
a = x - m(); // OK: dynamic binding (C++)
vector<int> v; // OK: template instantiation
```

Flow-of-Control Checks

```
myfunc()
{ ...
    break; // ERROR
}
```

```
myfunc()
{ ...
    while (n)
    { ...
        if (i>10)
            break; // OK
    }
}
```

```
myfunc()
{ ...
    switch (a)
    { case 0:
        ...
        break; // OK
        case 1:
        ...
    }
}
```

Uniqueness Checks

myfunc()
{ int i, j, i; // ERROR
...
}

```
cnufym(int a, int a) // ERROR
{ ...
```

```
struct myrec
{ int name;
};
struct myrec // ERROR
{ int id;
};
```

Outlines (Sections)

- 1. Variants of Syntax Trees
- 2. Three-Address Code
- 3. Types and Declarations
- 4. Translation of Expressions
- 5. Type Checking
- 6. Control Flow
- 7. Backpatching
- 8. Switch-Statements
- 9. Intermediate Code for Procedures

1. Variants of Syntax Trees

A directed acyclic graph (called a DAG) for an expression identifies the common subexpressions of the expression



2. Three-Address Code

In three-address code, there is at most one operator on the right side of an instruction. An address can be: name, constant, compiler-generated temporary.



$$t_{1} = b - c$$

$$t_{2} = a * t_{1}$$

$$t_{3} = a + t_{2}$$

$$t_{4} = t_{1} * d$$

$$t_{5} = t_{3} + t_{4}$$

Three-address code

Common Three-Address Instructions

- 1. Assignment instruction
- 2. Assignment
- 3. Copy instruction
- 4. Indexed copy instruction
- 5. Address and pointer assignment:
- 6. Unconditional jump
- 7. Conditional jump
- 9. Procedure call $p(x_1, x_2, \ldots, x_n)$: param x_1

```
x = y op z
                                    x = op y
                                    x = y
                                  x = y[i] and x[i] = y
                                  x = \&y, x = *y, and *x = y
                                    goto L
                                    if x relop y goto L
8. Conditional jump if x goto L and ifFalse x goto L
                                    param x_2
                                    ....
```

param x_n call p, n

Two Ways of Assigning Labels to Three-Address Statements



Quadruples, Triples, and Indirect Triples a = b*-c+b*-c;



3. Type Expressions

- A type expression is either a basic type or is formed by applying a type constructor to type expressions
 - Basic types: boolean, char, integer, float, etc.
 - Type constructors: pointer-to, array-of, records and classes, list-of, templates, and functions (s \rightarrow t).
 - Type names: typedefs in C and named types in Pascal
- Type expressions may contain variables whose values are type expressions



Cyclic Graph Representations

Source program

```
struct Node
{ int val;
   struct Node *next;
};
```

Internal compiler representation of the **Node** type: cyclic graph



Type Equivalence

- When type expressions are represented by graphs, two types are structurally equivalent if and only if one of the following conditions is true:
 - They are the same basic type.
 - They are formed by applying the same constructor to structurally equivalent types .
 - One is a type name that denotes the other.
- If type names are treated as standing for themselves, then the first two conditions in the above definition lead to name equivalence of type expressions

Structural Equivalence Example

- Two types are the same if they are *structurally identical*
- Used in C/C++, Java, C#



Type Equivalence Examples

struct Node
{ int val;
 struct Node *next;
};

struct Node s, *p; p = &s; // OK *p = s; // OK p = s; // ERROR



Storage Layout for Local Names

Type Declarations $\begin{array}{rcl} D & \rightarrow & T \text{ id } ; D & \mid \epsilon \\ T & \rightarrow & B C \mid \text{ record } '\{' D '\}' \\ B & \rightarrow & \text{ int } \mid \text{ float} \\ C & \rightarrow & \epsilon \mid [\text{ num }] C \end{array}$

record { int tag; float x; float y; } q;







Computing Types and Their Widths

Type Declarations

 \rightarrow T id; D | ϵ D $T \rightarrow BC \mid \mathbf{record} \ {}'{\{'D'\}'}$ $B \rightarrow \text{int} \mid \text{float}$ $C \rightarrow \epsilon \mid [\text{num}] C$

- $T \rightarrow B$
- $B \rightarrow int$
- $B \rightarrow \mathbf{float}$
- $C \rightarrow \epsilon$
- $\{ t = B.type; w = B.width; \}$ ${T.type = C.type; T.width = C.width; }$ $\{B.type = integer; B.width = 4;\}$ $\{B.type = float; B.width = 8; \}$ $\{C.type = t; C.width = w; \}$ $C \rightarrow [\text{num}] C_1 \quad \{ array(\text{num.value}, C_1.type); \}$ $C.width = \mathbf{num}.value \times C_1.width; \}$

Sequences of Declarations

Computing the relative addresses of declared names

$$P \rightarrow \{ top = new Evn(); offset = 0; \}$$

$$D \rightarrow T id ; \{ top.put(id.lexeme, T.type, offset);$$

$$offset = offset + T.width; \}$$

$$D \rightarrow \epsilon$$

Handling of field names in records

 $T \rightarrow \mathbf{record} '\{' \{ Env.push(top); top = \mathbf{new} Env(); \\ Stack.push(offset); offset = 0; \} \end{cases}$

 $D' \}' \{ T.type = record(top); T.width = offset;$ $top = Env.pop(); offset = Stack.pop(); \}$

Example: Annotated Parse Tree for int [2][3]



Example:

Determine types and relative addresses

float x; record { float x; float y; } p; record { int tag; float x; float y; } q;

D '}'	$\{ T.type = record(top); T.u$	width = offset;
	top = Env.pop(); offset =	Stack.pop(); }

 $Stack.push(offset): offset = 0; \}$

Line, id, type, offset, width

4. Translation of Expressions



 $a = t_2$



Figure 6.19: Three-address code for expressions

Translation of Expressions (cont.) Incremental Translation

$S \rightarrow id = E;$	{ gen(top.get(id.lexeme) '=' E.addr) ; }
$E \rightarrow E_1 + E_2$	{ E.addr = new Temp(); gen(E.addr '=' E ₁ .addr '+' E ₂ .addr) ; }
- E ₁	{ E.addr = new Temp() ; gen(E.addr '=' 'minus' E ₁ .addr) ; }
(E ₁)	{ E.addr = E_1 .addr; }
id	{ E.addr = top.get(id.lexeme) ;}

Figure 6.20: Generating three-address code for expressions incrementally

In the incremental approach, gen not only constructs a three-address instruction, but also appends the instruction to the sequence of instructions generated so far.

Addressing Array Elements

 $a[i].addr = base + i \times w$ $A[i_1][i_2].add$

 $\mathbf{A}[\mathbf{i}_1][\mathbf{i}_2].\mathbf{addr} = \mathbf{base} + \mathbf{i}_1 \times \mathbf{w}_1 + \mathbf{i}_2 \times \mathbf{w}_2$

 $A[i_1][i_2]...[i_k].addr = base + i_1 \times w_1 + i_2 \times w_2.... + i_k \times w_k$ (6.4)

Layouts for a two-dimensional array



Translation of Array References

 $S \rightarrow id = E$; { gen(top.get(id.lexeme) '=' E.addr); } L = E; { gen(L.addr.base'['L.addr']''='E.addr); } $E \rightarrow E_1 + E_2 \quad \{ E.addr = \mathbf{new} \ Temp(); \}$ $gen(E.addr'='E_1.addr'+'E_2.addr)$; } id $\{ E.addr = top.get(id.lexeme); \}$ $L \qquad \{ E.addr = \mathbf{new} \ Temp(); \}$ $gen(E.addr'='L.array.base'['L.addr']'); \}$ $L \rightarrow id [E] \{L.array = top.get(id.lexeme);$ L.type = L.array.type.elem; $L.addr = \mathbf{new} Temp();$ $gen(L.addr'='E.addr'*'L.type.width); \}$ $L_1 \ [E] \ \{L.array = L_1.array;$ $L.type = L_1.type.elem;$ $t = \mathbf{new} Temp();$ $L.addr = \mathbf{new} Temp();$ gen(t'='E.addr'*'L.type.width); $gen(L.addr'='L_1.addr'+'t);$

Figure 6.22: Semantic actions for array references

Translation of Array References (Cont.)

- Nonterminal L has three synthesized attributes:
 - 1. L.addr denotes a temporary that is used while computing the offset for the array reference by summing the $\mathbf{i_j} \times \mathbf{w_j}$ in (6.4)
 - 2. L.array is a pointer to the symbol-table entry for the array name.
 - L.array.base is the base address of the array.
 - L.array.type is the type of the array.
 - 3. L.type is the type of the subarray generated by L.
- Assume t is a type, then
 - t.width represents the width.
 - t.elem gives the element type.



5. Type Checking

- To do type checking a compiler needs to assign a type expression to each component of the source program.
- The compiler must then determine that these type expressions conform to a collection of logical rules that is called the type system for the source language
- Type checking can take on two forms:
 - Synthesis
 - Inference

Rules for Type Checking

- Type synthesis builds up the type of an expression from the types of its subexpressions.
- It requires names to be declared before they are used.

if f has type $s \rightarrow t$ and x has type s, then expression f (x) has type t

- Type inference determines the type of a language construct from the way it is used.
- It does not require names to be declared

if f(x) is an expression, then for some α and β , f has type $\alpha \rightarrow \beta$ and x has type α

Type Conversions

- Widening conversions
 - preserve information
- Narrowing conversions
 - lose information
- Coercions (implicit conversions)
 - are done automatically by the compiler.
- Casts (explicit conversions)
 - are done by programmer to write something to cause the conversion.





Introducing Type Conversions into Expression Evaluation

$$E \rightarrow E_1 + E_2 \quad \{ E.type = max(E_1.type, E_2.type); \\ a_1 = widen(E_1.addr, E_1.type, E.type); \\ a_2 = widen(E_2.addr, E_2.type, E.type); \\ E.addr = \mathbf{new} \ Temp(); \\ gen(E.addr'='a_1'+'a_2); \}$$

 $max(t_1, t_2)$ returns the maximum (or least upper bound) of the two types t_1 and t_2 in the widening hierarchy.

widen(a, t, w) generates type conversions
if needed to widen an address a of type t
into a value of type w.

x = 2 + 3.14

$$\downarrow$$

 $t_1 = (float) 2$
 $t_2 = t_1 + 3.14$
x = t_2

Overloading of Functions and Operators

Overloaded function examples

void err () { ... }
void err (String s) { ... }

A type-synthesis rule for overloaded functions

if *f* can have type $s_i \rightarrow t_i$, for $1 \le i \le n$, where $s_i \ne s_j$ for $i \ne j$ and *x* has type s_k , for some $1 \le k \le n$ then expression *f*(*x*) has type t_k

6.5.4 Type Inference and Polymorphic Functions

The term "polymorphic" refers to any code fragment that can be executed with arguments of different types

ML program for the length of a list

The type of *length*

$$\forall \alpha \ list(\alpha) \rightarrow integer$$

fun length(x) =**if** null(x) then 0 **else** length(tl(x)) + 1;

Example of use of *length*

length(["sun", "mon", "tue"]) + *length*([10, 9, 8, 7]) returns 7

Note: 6.5.4 and 6.5.5 need to be skipped for now.

Abstract syntax tree



Substitutions, Instances, and Unification

- A substitution *S* is a mapping from type variables to type expressions.
 - S(t) = the result of applying the substitution *S* to the variables in type expression *t*.
 - $S(\alpha) = integer$
 - $t = list(\alpha)$, then S(t) = list(integer)
 - $t = \alpha \rightarrow \alpha$, then $S(t) = integer \rightarrow integer$
- S(t) is called an instance of t.
- A substitution *S* is a *unifier* of two types t_1 and t_2 (t_1 and t_2 unify), if $S(t_1) = S(t_2)$.
- In the type inference algorithm, we *substitute* type variables by types to create type *instances*

Inferring a type for the function *length*

fun length(x) = if null(x) then 0 else length(tl(x)) + 1;

LINE	EXPRESSION	:	TYPE	UNIFY
1)	length	:	$\beta \rightarrow \gamma$	
2)	x	:	β	
3)	if	:	$boolean \times \alpha_i \times \alpha_i \to \alpha_i$	
4)	null	:	$list(\alpha_n) \rightarrow boolean$	
5)	null(x)	:	boolean	$list(\alpha_n) = \beta$
6)	0	:	integer	$\alpha_i = integer$
7)	+	:	$integer \times integer \rightarrow integer$	
8)	tl	:	$list(\alpha_t) \rightarrow list(\alpha_t)$	
9)	tl(x)	:	$list(\alpha_t)$	$list(\alpha_t) = list(\alpha_n)$
10)	length(tl(x))	:	γ	$\gamma = integer$
11)	1	:	integer	
12)	length(tl(x)) + 1	:	integer	
13)	if (· · ·)	:	integer	

 $\forall \alpha_n. \text{ list}(\alpha_n) \rightarrow integer$

6.5.5: An Algorithm for Unification

Examples 6.18: Consider the two type Expressions t_1 , t_2 , and the substitution *S*

 $\begin{array}{l} \mathsf{t}_1 = ((\alpha_1 \to \alpha_2) \times \mathit{list}(\alpha_3)) \to \mathit{list}(\alpha_2) \\ \mathsf{t}_2 = ((\alpha_3 \to \alpha_4) \times \mathit{list}(\alpha_3)) \to \alpha_5 \end{array}$

$$\begin{array}{c|c} x & S(x) \\ \hline \alpha_1 & \alpha_1 \\ \alpha_2 & \alpha_2 \\ \alpha_3 & \alpha_1 \\ \alpha_4 & \alpha_2 \\ \alpha_5 & list(\alpha_2) \end{array}$$

$$S(t_1) = S(t_2) = ((\alpha_1 \to \alpha_2) \times list(\alpha_1)) \to list(\alpha_2)$$



An Algorithm for Unification(Cont.)

```
boolean unify(Node \ m, Node \ n) {
      s = find(m); t = find(n);
      if (s = t) return true;
      else if ( nodes s and t represent the same basic type ) return true;
       else if (s is an op-node with children s_1 and s_2 and
                    t is an op-node with children t_1 and t_2 {
              union(s,t);
              return unify(s_1, t_1) and unify(s_2, t_2);
       }
       else if s or t represents a variable {
              union(s,t);
              return true;
       }
       else return false;
```