


复旦大学

编译原理/Compiler Semantic Analysis


周雅倩
复旦大学计算机科学技术学院
zhouyaqian@fudan.edu.cn
2018/11/9



References

- <http://www.stanford.edu/class/cs143/>
 - Partial contents are copied from the slides of Alex Aiken


2



Outline

- Symbol Table
- Type Check


2018/11/9 3



Abstract Syntax Trees

- Separate AST construction from semantic checking phase
- Traverse the AST and perform semantic checks (or other actions) only after the tree has been built and its structure is stable
- This approach is less error-prone
 - It is better when efficiency is not a critical issue


2018/11/9 4



Incorrect Programs

- Lexically and syntactically correct programs may still contain other errors!
- Lexical and syntax analysis are not powerful enough to ensure the correct usage of variables, objects, functions, statements, etc.

2018/11/9 5



Incorrect Programs

- Example 1: lexical analysis does not distinguish between different variable or function identifiers (it returns the same token for all identifiers)


```
int a;      int a;
a = 1;      b = 1;
```
- Example 2: syntax analysis does not correlate the declarations with the uses of variables in the program:


```
int a;      a = 1;
a = 1;
```
- Example 3: syntax analysis does not correlate the types from the declarations with the uses of variables:


```
int a;      int a;
a = 1;      a = 1.0;
```

2018/11/9 6

Goals of Semantic Analysis

- **Semantic analysis** = ensure that the program satisfies a set of rules regarding the usage of programming constructs (variables, objects, expressions, statements)
- Examples of semantic rules:
 - Variables must be declared before being used
 - A variable should not be declared multiple times in the same scope
 - In an assignment statement, the variable and the assigned expression must have the same type
 - The condition of an if-statement must have type boolean
- Some categories of rules:
 - Semantic rules regarding types
 - Semantic rules regarding scopes

2018/11/9 7

Type Information

- **Type information** = describes what kind of values correspond to different constructs: variables, statements, expressions, functions
- variables: `int a;` **integer**
- expressions: `(a+1) == 2` **boolean**
- statements: `a = 1.0` **floating-point**
- functions: `int pow(int n, int m)` **`int x int → int`**

2018/11/9 8

Type Checking (cont.)

- **Type checking** = set of rules that ensure the type consistency of different constructs in the program
- Examples:
 - The type of a **variable** must match the type from its declaration
 - The operands of **arithmetic expressions** (+, *, -, /) must have integer types; the result has integer type
 - The operands of **comparison expressions** (==, !=) must have integer or string types; the result has boolean type

2018/11/9 9

Type Checking (cont.)

- More examples:
 - For each **assignment statement**, the type of the updated variable must match the type of the expression being assigned
 - For each **call statement** `foo(v1, ..., vn)`,
 - the type of each actual argument **v_i** must match the type of the corresponding formal argument **f_i** from the declaration of function **foo**
 - the type of the return value must match the return type from the declaration of the function

2018/11/9 10

Scope Information

- **Scope information** = characterizes the declaration of identifiers and the portions of the program where it is allowed to use each identifier
 - Example identifiers: variables, functions, objects, labels
- **Lexical scope** = textual region in the program
 - Statement block
 - Formal argument list
 - Object body
 - Function or method body
 - Module body
 - Whole program (multiple modules)
- **Scope of an identifier**: the lexical scope in which it is valid

2018/11/9 11

Scope of variables in statement blocks

```

{ int a;
  ...
  { int b; } b
  ...
}

```

Diagram illustrating scope: The inner block `{ int b; }` is enclosed within the outer block `{ int a; ... }`. The variable `b` is shown with a bracket indicating its scope is limited to the inner block, while `a` is shown with a bracket indicating its scope covers the entire outer block.

- In C:
 - Scope of file static variables: current file
 - Scope of external variables: whole program
 - Scope of automatic variables, formal parameters, and function static variables: the function

2018/11/9 12

Scope of formal arguments of functions/methods

```

int factorial(int n) {
    ...
}

```

Scope of **labels**

```

void f() {
    ... goto l; ...
l: a = 1;
    ... goto l; ...
}

```

2018/11/9 13

Scope of object fields and methods

```

class A {
    private int x;
    public void g() { x = 1; }
    ...
}
class B extends A {
    ...
    public int h() { g(); }
    ...
}

```

2018/11/9 14

Semantic Rules for Scopes

- Main rules regarding scopes:
 - Rule 1: Use an identifier only if defined in enclosing scope
 - Rule 2: Do not declare identifiers of the same kind with identical names more than once in the same lexical scope
- Can declare identifiers with the same name with identical or overlapping lexical scopes if they are of different kinds

```

class X {
    int X;
    void X(int X) {
        X: for(;;)
            break X;
    }
}
int X(int X) {
    int X;
    goto X;
    { int X;
      X: X = 1; }
}

```

2018/11/9 15

Not Recommended!

SYMBOL TABLE

2018/11/9 16

Symbol Tables

- Semantic checks refer to properties of identifiers in the program -- their scope or type
- Need an environment to store the information about identifiers = **symbol table**
- Each entry in the symbol table contains
 - the name of an identifier
 - additional information: its kind, its type, if it is constant, ...

2018/11/9 17

Symbol Tables(Example)

| NAME | KIND | TYPE | ATTRIBUTES |
|------|------|------------------|------------|
| foo | fun | int * int → bool | extern |
| m | arg | int | |
| n | arg | int | const |
| tmp | var | bool | const |

2018/11/9 18

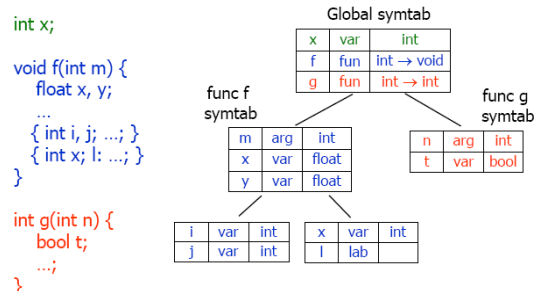
Scope Information

- How to capture the scope information in the symbol table?
- Idea:
 - There is a hierarchy of scopes in the program
 - Use a similar **hierarchy of symbol tables**
 - One symbol table for each scope
 - Each symbol table contains the symbols declared in that lexical scope

2018/11/9

19

Example



2018/11/9

20

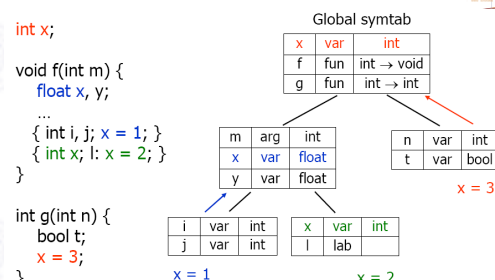
Identifiers with Same Name

- The hierarchical structure of symbol tables automatically solves the problem of resolving name collisions (identifiers with the same name and overlapping scopes)
- To find which is the declaration of an identifier that is active at a program point:
 - Start from the current scope
 - Go up in the hierarchy until you find an identifier with the same name

2018/11/9

21

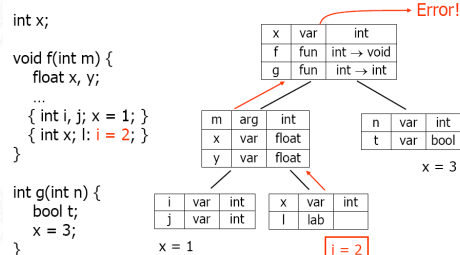
Example



2018/11/9

22

Catching Semantic Errors



2018/11/9

23

Symbol Table Operations

- Two operations:
 - insert**
 - lookup**
- Each entry can be implemented as a record. Records can have different formats (Variant records in Pascal).

编译原理

Symbol Table Operations

- Cannot build symbol tables during lexical analysis
 - hierarchy of scopes encoded in the syntax
- Build the symbol tables:
 - while parsing, using the semantic actions
 - After the AST is constructed

25

编译原理

Storing characters

- Method 1: A fixed size space within each entry large enough to hold the largest possible name.
 - Most names will be much shorter than this so there will be a lot of wasted storage
- Method 2: Store all symbols in one large separate array. Each symbol is terminated with an end of symbol mark (EOS).
 - Each symbol table record contains a pointer to the first character of the symbol.
- Method n: modern languages (e.g. Java, C++ std components) has efficient DS, e.g. string or vector

编译原理

Symbol Table Data Structure

- One Linear list:
 - Easy to implement
 - Search time will be very long if source has many symbols.
- Simple implementation = array
 - One entry per symbol
 - Disadvantage:
 - table has fixed size
 - need to know in advance the number of entries
- Dynamic structure = list
 - Can grow dynamically during compilation
 - need to scan half the list on average

编译原理

Symbol Table Data Structure

Hash table:

- Run the symbol name through a hash function to create an index in a table.
- If some other symbol has already claimed the space then rehash with another hash function to get another index, etc.
- Hash Table must be large enough to accommodate largest number of symbols.

编译原理

Symbol Table Data Structure

- Open hash:
 - Store the entries in a number of linear lists (called **Buckets**).
 - Use a hash function on the symbol name to determine which lists to use.
 - A good hash function will spread the symbols across the buckets, so each linear list will be short.

编译原理

Forward References

- Forward references = use an identifier within the scope of its declaration, but before it is declared
 - Forward declaration
- Any compiler phase that uses the information from the symbol table must be performed after the table is constructed
 - Cannot type-check and build symbol table at the same time

Example:

```
class A {
    int m() { return n(); }
    int n() { return 1; }
}
```

2018/11/9

30

编译原理

Summary

- **Semantic checks** ensure the correct usage of variables, objects, expressions, statements, functions, and labels in the program
- **Scope semantic checks** ensure that identifiers are correctly used within the scope of their declaration
- **Type semantic checks** ensures the type consistency of various constructs in the program
- **Symbol tables**: a data structure for storing information about symbols in the program
 - Used in semantic analysis and subsequent compiler stages

2018/11/9 31

编译原理

TYPE CHECKING

2018/11/9 32

编译原理

What Are Types?

- **Types** describe the values computed during the execution of the program
- Essentially, types are predicate on values, e.g., “int x” in Java means “ $x \in [-2^{31}, 2^{31})$ ”
- **Type errors**: improper, type-inconsistent operations during program execution
- **Type-safety**: absence of type errors

2018/11/9 33

编译原理

How to Ensure Type-Safety

- Bind (assign) types, then check types
- **Type binding**: defines type of constructs in the program (e.g., variables, functions)
 - Can be either explicit (int x) or implicit ($x = 1$)
 - Type consistency (safety) = correctness with respect to the type bindings
- **Type checking**: determine if the program correctly uses the type bindings
 - Consists of a set of type-checking rules

2018/11/9 34

编译原理

Type Checking

- Type checking: static semantic checks to enforce the type safety of the program
- Examples:
 - Unary and binary operators (e.g., +, ==, []) must receive operands of the proper type
 - Functions must be invoked with the right number and type of arguments
 - Return statements must agree with the return type
 - In assignments, assigned value must be compatible with type of variable on LHS.
 - Class members accessed appropriately

2018/11/9 35

编译原理

Static vs. Dynamic Typing

- Static and dynamic typing refer to type definitions (i.e., bindings of types to variables, expressions, etc.)
 - Statically typed language: types are defined and checked at compile-time, and do not change during the execution of the program
 - E.g., C, Java, Pascal
 - Dynamically typed language: types defined and checked at run-time, during program execution
 - E.g., Lisp, Smalltalk

2018/11/9 36

Strong vs. Weak Typing

- Strong and weak typing refer to how much type consistency is enforced
 - Strongly typed languages:** guarantees that accepted programs are type-safe
 - Weakly typed languages:** allow programs that contain type errors
- Can achieve strong typing using either static or dynamic typing

2018/11/9 37

Soundness

- Soundness type systems:** can statically ensure that the program is type-safe.
- Soundness implies strong typing.
- Static type safety requires a conservative approximation of the values that may occur during all possible executions
 - May reject type-safe programs
 - Need to be expressive: reject as few type-safe programs as possible

2018/11/9 38

Why Static Checking?

- Efficient code
 - Dynamic checks slow down the program
- Guarantees that all executions will be safe
 - Dynamic checking gives safety guarantees only for some execution of the program
- But is conservative for sound systems
 - Needs to be expressive: reject few type-safe programs

2018/11/9 39

Type Systems

- Type is predicate on value
 - Type expressions:** describe the possible types in the program: int, string, array[], Object, etc.
 - Type system:** defines types for language constructs (e.g., expressions, statements)

2018/11/9 40

Type Systems

- The design of a type checker for a language is based on
 - the syntactic constructs in the language
 - the notion of types
 - the rules of assigning types to language constructs
- In C and Pascal
 - arithmetic
 - pointer operator
- Each expression have a type associated it.
- Types have structure.

2018/11/9 41

Kinds of Types

- Basic Type:**
 - boolean, char, int, real, enum
- Constructed Type:**
 - arrays, records, sets, pointers, functions

2018/11/9 42

Type Expressions

- Basic type
- Type name
- Type Constructor

2018/11/9 43

Basic type

- boolean, char, int, real, enum
- void
 - the absence of a value
 - allows statements to be checked
- type_error
 - Special basic type
 - Will signal an error during type checking

2018/11/9 44

Type Constructors

- Can be applied to type expressions so that a new type expression generated
- Typical type constructors
 - Arrays: `array(I, T)`
 - Products: `T1 x T2`
 - Structures: `record((id1 x T1) x (id2 x T2) x ...)`
 - Pointers: `pointer(T)`
 - Functions: `T1 x ... x Tn → T`

2018/11/9 45

Type Constructors

- Various kinds of array types in different
 - `array(T, S)`
 - T : element type
 - S : index set, often a range of integers
- programming languages
 - `array(T)` : arrays without bounds
 - C, Java: `T [], Module-3: array of T`
 - `array(T, S)` : array with size
 - C: `T[S]`, Module-3: `array[S] of T`
 - May be indexed 0..S-1
 - `array(T, L, U)` : array with upper/lower bounds
 - Pascal: `array[L..U] of T`
 - `array(T, S1, ..., Sn)` : multi-dimensional arrays
 - FORTRAN: `T(L1..., Ln)`

2018/11/9 46

Structures(Records)

Structures in C:

```
struct { int a; float b; }
```

Records in Pascal:

```
record a: integer; b: real;
end
```

```
struct list {
  struct list* next;
  int value;
}
```

- Field names should be part of the type constructor
- Sometime it is more convenient to represent record as the following:
 - `record((id1 x T1) x (id2 x T2) x ...)`

2018/11/9 47

Type Constructor

- Pointers:
 - `pointer(T)`
 - pointer to an object of type T
- Functions:
 - Mapping elements of domain to range
 - `T1 x ... x Tn → T`

C pointers:

```
T* (e.g., int *x);
```

Pascal pointers:

```
^T (e.g., x: ^integer);
```

C functions:

```
int f(float x, float y)
```

2018/11/9 48

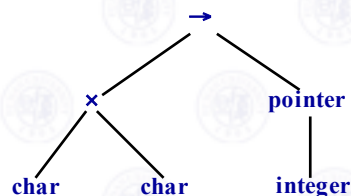
Type names

- Type expressions may be named
- A type name is a type expression

2018/11/9

49

Type Expressions



2018/11/9

50

Type Systems

- A Type System is
 - a collection of rules for
 - assigning type expressions to
 - the various parts of a program

2018/11/9

51

Type-Checking

- (Option) first build the AST, then implement type-checking by recursive traversal of the AST nodes:

```

class Add extends Expr {
  Type typeCheck (Symtab s) {
    Type t1 = e1.typeCheck(s);
    t2 = e2.typeCheck(s);
    if (t1 == Int && t2 == Int) return Int;
    else throw new TypeCheckError("+");
  }
}
  
```

2018/11/9

52

Static and Dynamic Checking of Types

- Static:
 - Checking done by a compiler
- Dynamic:
 - Checking done when the target program runs
- A sound type system eliminates the need for dynamic checking for type errors.

2018/11/9

53

Error Recovery

- The compiler must report the nature and location of error.
- It is desirable for the type checker to recover from errors.

2018/11/9

54

A Simple Language

- $P \rightarrow D ; E$
- $D \rightarrow D ; D \mid id : T$
- $T \rightarrow char \mid integer \mid array[num] \text{ of } T \mid *T$
- $E \rightarrow literal \mid num \mid id \mid E \text{ mod } E \mid E[E] \mid *E$
- **key: integer;**
- **key mod 1999**

2018/11/9 55

Semantic Actions

- **id : char**
- **addtype(id.entry, char)**
- **id : integer**
- **addtype(id.entry, integer)**

2018/11/9 56

Semantic Actions

- $T \rightarrow array[num] \text{ of } T1$
- **T.type = array(1..num.val, T1.type)**
- $T \rightarrow *T1$
- **T.type = pointer(T1.type)**

2018/11/9 57

Type Checking of Expressions

- $E \rightarrow literal$
- **E.type = char**
- $E \rightarrow num$
- **E.type = integer**
- $E \rightarrow id$
- **E.type = lookup(id.entry)**

2018/11/9 58

Type Checking of Expressions

- $E \rightarrow E1 \text{ mod } E2$
- **E.type = (E1.type == integer && E2.type == integer) ? integer : type_error**

2018/11/9 59

Type Checking of Expressions

- $E \rightarrow E1[E2]$
- **E.type = (E2.type == integer && E1.type == array(s, t)) ? t : type_error**

2018/11/9 60

Type Checking of Expressions

- $E \rightarrow *E1$
- $E.type = (E1.type == \text{pointer}(t)) ? t : \text{type_error}$

2018/11/9 61

Type Checking of Functions

- $T \rightarrow T1 \rightarrow T2$
- $T.type = T1.type \rightarrow T2.type$
- $E \rightarrow E1(E2)$
- $E.type = (E2.type == s \ \&\& \ E1.type == s \rightarrow t) ? t : \text{type_error}$

2018/11/9 62

Names for Type Expressions

- `typedef cell *link;`
- `link next;`
- `link last;`
- `cell *p;`
- `cell *q, *r;`
- `next, last` `link`
- `p, q, r` `pointer(cell)`

2018/11/9 63

Names for Type Expressions

- **Name Equivalence**
 - Views each type name as a distinct type
- **Structural Equivalence**
 - Names are replaced by the type expressions they defined

2018/11/9 64

Structural Equivalence of Type Expressions

- Two expressions are
 - Either the same basic type
 - Or are formed by applying the same constructor to structurally equivalent types

2018/11/9 65

Algorithm

```

int sequiv(s, t) {
    if ( s and t are the same basic type ) return 1 ;
    else if ( s == array(s1, s2) && t == array(t1, t2) )
        return sequiv(s1, t1) && sequiv(s2, t2)
    else if ( s == s1xs2 && t == t1xt2 )
        return sequiv(s1, t1) && sequiv(s2, t2)
    else if ( s == pointer(s1) && t == pointer(t1) )
        return sequiv(s1, t1)
    else if ( s == s1→s2 && t == t1→t2 )
        return sequiv(s1, t1) && sequiv(s2, t2)
    return 0;
}

```

2018/11/9 66

Names for Type Expressions

- `typedef cell *link;`
- `link next;`
- `link last;`
- `cell *p;`
- `cell *q, *r;`

■ `next, last` `link`
 ■ `p, q, r` `pointer(cell)`

2018/11/9 67

Type Variables

- Variables representing type expressions allow us to talk about unknown types.

2018/11/9 68

Names for Type Expressions

2018/11/9 69

Cycles in Representations of Types

```
typedef cell *link ;
typedef struct {
    int info ;
    link next ;
} cell ;
```

2018/11/9 70

Cycles in Representations of Types

2018/11/9 71

Cycles in Representations of Types

2018/11/9 72

Type Conversion

- **X + I**
 - X real, I int
 - Insert into real before I
 - X I into real real+
- **Implicit type conversion**
 - **Coercions**
 - No information lost
 - int to double, not vice-versa
- **Explicit type casting**

73

Type Conversion

| Production | Semantic rule |
|---------------------------------------|---|
| $E \rightarrow \text{num}$ | $E.\text{type} = \text{integer}$ |
| $E \rightarrow \text{num}.\text{num}$ | $E.\text{type} = \text{real}$ |
| $E \rightarrow \text{id}$ | $E.\text{type} = \text{lookup}(\text{id}.\text{entry})$ |
| $E \rightarrow E_1 \text{ op } E_2$ | $E.\text{type} = \text{if}(E_1.\text{type} = \text{integer} \ \&\& \ E_2.\text{type} = \text{integer}) \text{ then integer}$ $\quad \text{else if}(E_1.\text{type} = \text{integer} \ \&\& \ E_2.\text{type} = \text{real}) \text{ then real}$ $\quad \text{else if}(E_1.\text{type} = \text{real} \ \&\& \ E_2.\text{type} = \text{integer}) \text{ then real}$ $\quad \text{else if}(E_1.\text{type} = \text{real} \ \&\& \ E_2.\text{type} = \text{real}) \text{ then real}$ $\quad \text{else type_error}$ |

74

Simple Overloading

- **Overloaded symbol has different meanings depending on its arguments**
 - 1.0+2.0
 - 1 + 2
- **Resolution**
 - Determine a unique meaning for an occurrence of an overloading symbol
 - Well known as operator identification

75

Narrowing the Set of Possible Types

- **Given a unique type from the context, we narrow down the type choices for each subexpression.**

76

General Overloading

- The meaning of an overloading symbol is depending on its context
- An example from Ada, operator “*”
 - A pair of integers to an integer
 - A pair of integers to a complex
 - A pair of complexes to a complex
 - $2^*(3*5)$
 - $(3*5)$ must be integer
 - $(3*5)*z$
 - $(3*5)$ must be complex

77

General Overloading

- $E \rightarrow E_1(E_2)$
- $E.\text{type} = \{(E_2.\text{types} == s \ \&\& \ E_1.\text{types} == s \rightarrow t) ? t : \text{type_error}\}$
- $E.\text{types} = \{ t \mid \text{there exists an } s \text{ in } E_2.\text{types} \text{ such that } (s \rightarrow t) \text{ is in } E_1.\text{types} \}$

78

