

Block STRUCTURE CONCEPTS

- Nested Visibility
No access to identifiers outside their scope.
- Nearest Declaration Applies
Using static nesting of scopes.
- Automatic Allocation and Deallocation of Locals
Lifetime of data objects is bound to the scope of the Identifiers that denote them.

Block-STRUCTURED Symbol Tables

Block structured symbol tables are designed to support the scoping rules of block structured languages.

For our CSX project we'll use class **Symb** to represent symbols and **SymbolTable** to implement operations needed for a block-structured symbol table.

Class **Symb** will contain a method

```
public String name()
```

that returns the name associated with a symbol.

Class **SymbolTable** contains the following methods:

- **public void openScope()** {
A new and empty scope is opened.
- **public void closeScope() throws EmptySTEException**
The innermost scope is closed. An exception is thrown if there is no scope to close.
- **public void insert(Symb s) throws DuplicateException, EmptySTEException**
A **Symb** is inserted in the innermost scope. An exception is thrown if a **Symb** with the same name is already in the innermost scope or if there is no symbol table to insert into.

- **public Symb localLookup(String s)**

The innermost scope is searched for a **Symb** whose name is equal to **s**. Null is returned if no **Symb** named **s** is found.

- **public Symb globalLookup(String s)**

All scopes, from innermost to outermost, are searched for a **Symb** whose name is equal to **s**. The first **Symb** that matches **s** is found; otherwise null is returned if no matching **Symb** is found.

Is CASE SIGNIFICANT?

In some languages (C, C++, Java and many others) case *is* significant in identifiers. This means `aa` and `AA` are different symbols that may have entirely different definitions.

In other languages (Pascal, Ada, Scheme, CSX) case *is not* significant. In such languages `aa` and `AA` are two alternative spellings of the same identifier.

Data structures commonly used to implement symbol tables usually treat different cases as different symbols. This is fine when case is significant in a language. When case is insignificant, you probably will

need to *strip case* before entering or looking up identifiers.

This just means that identifiers are converted to a uniform case before they are entered or looked up. Thus if we choose to use lower case uniformly, the identifiers `aaa`, `AAA`, and `AaA` are all converted to `aaa` for purposes of insertion or lookup.

BUT, inside the symbol table the identifier is stored in the form it was declared so that programmers see the form of identifier they expect in listings, error messages, etc.

How ARE Symbol Tables IMPLEMENTED?

There are a number of data structures that can reasonably be used to implement a symbol table:

- An Ordered List
Symbols are stored in a linked list, sorted by the symbol's name. This is simple, but may be a bit too slow if many identifiers appear in a scope.
- A Binary Search Tree
Lookup is much faster than in linked lists, but rebalancing may be needed. (Entering identifiers in sorted order turns a search tree into a linked list.)
- Hash Tables
The most popular choice.

IMPLEMENTING Block-STRUCTURED Symbol Tables

To implement a block structured symbol table we need to be able to efficiently open and close individual scopes, and limit insertion to the innermost current scope.

This can be done using one symbol table structure if we tag individual entries with a "scope number."

It is far easier (but more wasteful of space) to allocate one symbol table for each scope. Open scopes are stacked, pushing and popping tables as scopes are opened and closed.

Be careful though—many preprogrammed stack implementations don't allow you to “peek” at entries below the stack top. This is necessary to lookup an identifier in all open scopes.

If a suitable stack implementation (with a peek operation) isn't available, a linked list of symbol tables will suffice.

n is the length of the string, c_i is the i -th character and all arithmetic is done without overflow checking.

Why such an elaborate hash function?

Simpler hash functions can have major problems.

Consider $\sum_{i=0}^{n-1} c_i$ (add the characters).

For short identifiers the sum grows slowly, so large indices won't often be used (leading to non-uniform use of the hash table).

MORE ON HASHTABLES

Hashtables are a very useful data structure. Java provides a predefined `Hashtable` class. Python includes a built-in `dictionary` type.

Every Java class has a `hashCode` method, which allows any object to be entered into a Java `Hashtable`.

For most objects, hash codes are pretty simple (the address of the corresponding object is often used).

But for strings Java uses a much more elaborate hash

$$\text{function: } \sum_{i=0}^{n-1} c_i \times 37^i$$

We can try $\prod_{i=0}^{n-1} c_i$ (product of characters), but now (surprisingly) the size of the hash table becomes an issue. The problem is that if even one character is encoded as an even number, the product *must* be even.

If the hash table is even in size (a natural thing to do), most hash table entries will be at even positions. Similarly, if even one character is encoded as a multiple of 3, the whole product will be a multiple of 3, so hash tables that are a multiple of three in size won't be uniformly used.

To see how bad things can get, consider a hash table with size 210 (which is equal to $2 \times 3 \times 5 \times 7$). This should be a particularly bad table size if a product hash is used. (Why?)

Is it? As an experiment, all the words in the Unix spell checker's dictionary (26000 words) were entered. Over 50% (56.7% actually) hit position 0 in the table!

Why such non-uniformity?

If an identifier contains characters that are multiples of 2, 3, 5 and 7, then their hash will be a multiple of 210 and will map to position 0.

For example, in `wisconsin`, `n` has an ASCII code of 110 (2×55) and `i` has a code of 105 (7×3).

If we change the table size ever so slightly, to 211, no table entry gets more than 1% of the 26000 words hashed, which is very good.

Why such a big difference? Well 211 is *prime* and there is a bit a folk-wisdom that states that prime numbers are good choices for hash table sizes. Now our product hash will cover table entries far more uniformly (small factors in the hash don't divide the table size evenly).

Now the reason for Java's more complex string hash function becomes evident—it can uniformly fill a hash table whose size isn't prime.

How ARE Collisions HANDLED?

Since identifiers are often unlimited in length, the set of possible identifiers is infinite. Even if we limit ourselves to short identifiers (say 10 or fewer characters), the range of valid identifiers is greater than 26^{10} .

This means that all hash tables need to contend with *collisions*, when two different identifiers map to the same place in the table.

How are collisions handled?

The simplest approach is *linear resolution*. If identifier `ia` hashes to position `p` in a hash

table of size s and position p is already filled, we try $(p+1) \bmod s$, then $(p+2) \bmod s$, until a free position is found.

As long as the table is not too filled, this approach works well. When we approach an almost-filled situation, long search chains form, and we degenerate to an unordered list.

If the table is 100% filled, linear resolution fails.

Some hash table implementations, including Java's, set a *load factor* between 0 and 1.0. When the fraction of filled entries in the table exceeds the load factor,

table size is increased and all entries are rehashed.

Note that bundling of a `hashCode` method within all Java objects makes rehashing easy to do automatically. If the hash function is external to the symbol table entries, rehashing may need to be done manually by the user.

An alternative to linear resolution is *chained resolution*, in which symbol table entries contain pointers to chains of symbols rather than a single symbol. All identifiers that hash to the same position appear on the same chain. Now overflowing table size is not catastrophic—

as the table fills, chains from each table position get longer. As long as the table is not too overfilled, average chain length will be small.

READING ASSIGNMENT

Read Chapter 3 of
Crafting a Compiler Second Edition.

SCANNING

A scanner transforms a character stream into a token stream. A scanner is sometimes called a *lexical analyzer* or *lexer*. Scanners use a formal notation (*regular expressions*) to specify the precise structure of tokens. But why bother? Aren't tokens very simple in structure? Token structure can be more detailed and subtle than one might expect. Consider simple quoted strings in C, C++ or Java. The body of a string can be any sequence of characters *except* a quote character (which must be escaped). But is this simple definition really correct?

Can a newline character appear in a string? In C it cannot, unless it is escaped with a backslash.

C, C++ and Java allow escaped newlines in strings, Pascal forbids them entirely. Ada forbids *all* unprintable characters.

Are null strings (zero-length) allowed? In C, C++, Java and Ada they are, but Pascal forbids them.

(In Pascal a string is a packed array of characters, and zero length arrays are disallowed.)

A precise definition of tokens can ensure that lexical rules are clearly stated and properly enforced.

REGULAR EXPRESSIONS

Regular expressions specify simple (possibly infinite) sets of strings. Regular expressions routinely specify the tokens used in programming languages.

Regular expressions can drive a *scanner generator*.

Regular expressions are widely used in computer utilities:

- The Unix utility *grep* uses regular expressions to define search patterns in files.
- Unix shells allow regular expressions in file lists for a command.

- Most editors provide a “context search” command that specifies desired matches using regular expressions.

- The Windows Find utility allows some regular expressions.

REGULAR SETS

The sets of strings defined by *regular expressions* are called *regular sets*.

When scanning, a token class will be a regular set, whose structure is defined by a regular expression.

Particular instances of a token class are sometimes called *lexemes*, though we will simply call a string in a token class an *instance* of that token. Thus we call the string `abc` an identifier if it matches the regular expression that defines valid identifier tokens.

Regular expressions use a finite character set, or *vocabulary* (denoted Σ).

This vocabulary is normally the character set used by a computer. Today, the *ASCII* character set, which contains a total of 128 characters, is very widely used.

Java uses the *Unicode* character set which includes all the ASCII characters as well as a wide variety of other characters.

An empty or *null* string is allowed (denoted λ , “lambda”). Lambda represents an empty buffer in which no characters have yet been matched. It also represents optional parts of tokens. An integer literal may begin with a plus or minus, or it may begin with λ if it is unsigned.