

CSE 309 (Compilers)

Basic Concepts

Dr. Muhammad Masroor Ali

Professor

Department of Computer Science and Engineering
Bangladesh University of Engineering and Technology
Dhaka-1000, Bangladesh

July, 2008

Version: 1.3, Last modified: August 3, 2008

Introduction

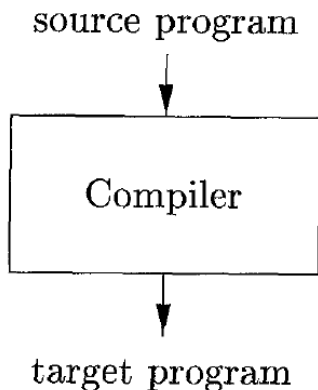
- Programming languages are notations for describing computations to people and to machines.
- The world as we know it depends on programming languages, because all the software running on all the computers was written in some programming language.
- But, before a program can be run, it first must be translated into a form in which it can be executed by a computer.
- The software systems that do this translation are called compilers.

Introduction

- Programming languages are notations for describing computations to people and to machines.
- The world as we know it depends on programming languages, because all the software running on all the computers was written in some programming language.
- But, before a program can be run, it first must be translated into a form in which it can be executed by a computer.
- The software systems that do this translation are called compilers.

Language Processors

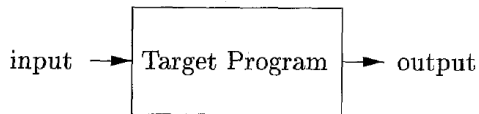
- Simply stated, a compiler is a program that can read a program in one language — the source language — and translate it into an equivalent program in another language - the target language.
- An important role of the compiler is to report any errors in the source program that it detects during the translation process.



A compiler

Language Processors

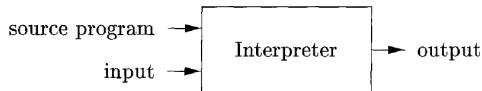
- If the target program is an executable machine-language program, it can then be called by the user to process inputs and produce outputs.



Running the target program

Language Processors — *continued*

- An interpreter is another common kind of language processor.
- Instead of producing a target program as a translation, an interpreter appears to directly execute the operations specified in the source program on inputs supplied by the user.

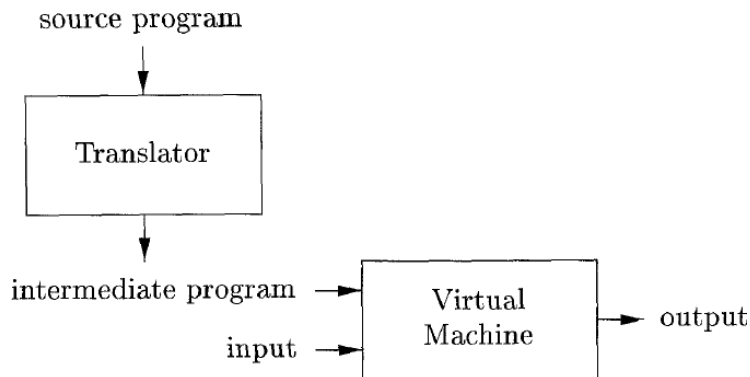


An interpreter

Language Processors — *continued*

- The machine-language target program produced by a compiler is usually much faster than an interpreter at mapping inputs to outputs.
- An interpreter, however, can usually give better error diagnostics than a compiler, because it executes the source program statement by statement.

Example



A hybrid compiler

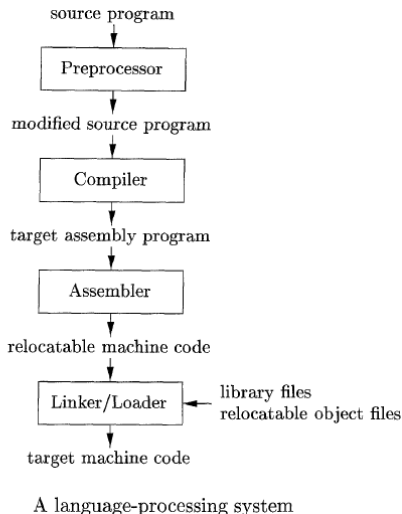
- Java language processors combine compilation and interpretation.

Example — *continued*

- A Java source program may first be compiled into an intermediate form called bytecodes.
- The bytecodes are then interpreted by a virtual machine.
- A benefit of this arrangement is that bytecodes compiled on one machine can be interpreted on another machine, perhaps across a network.
- In order to achieve faster processing of inputs to outputs, some Java compilers, called **just-in-time** compilers, translate the bytecodes into machine language immediately before they run the intermediate program to process the input.

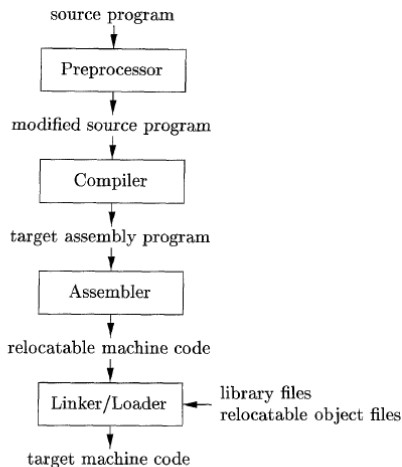
Language Processors — *continued*

- In addition to a compiler, several other programs may be required to create an executable target program.
- A source program may be divided into modules stored in separate files.
- The task of collecting the source program is sometimes entrusted to a separate program, called a preprocessor.
- The preprocessor may also expand shorthands, called macros, into source language statements.



Language Processors — *continued*

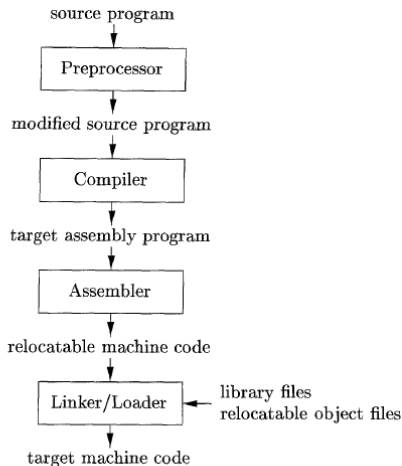
- The modified source program is then fed to a compiler.
- The compiler may produce an assembly-language program as its output, because assembly language is easier to produce as output and is easier to debug.
- The assembly language is then processed by a program called an assembler that produces relocatable machine code as its output.



A language-processing system

Language Processors — *continued*

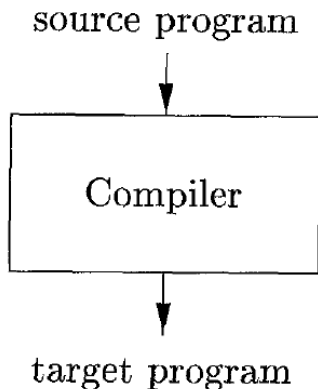
- Large programs are often compiled in pieces, so the relocatable machine code may have to be linked together with other relocatable object files and library files into the code that actually runs on the machine.
- The linker resolves external memory addresses, where the code in one file may refer to a location in another file.
- The loader then puts together all of the executable object files into memory for execution.



A language-processing system

The Structure of a Compiler

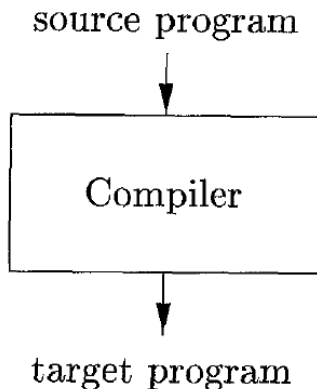
- Up to this point we have treated a compiler as a single box that maps a source program into a semantically equivalent target program.
- If we open up this box a little, we see that there are two parts to this mapping: analysis and synthesis.



A compiler

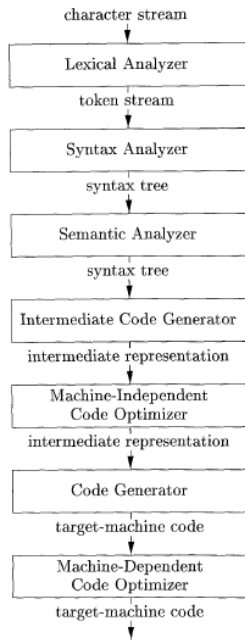
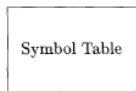
The Structure of a Compiler

- Up to this point we have treated a compiler as a single box that maps a source program into a semantically equivalent target program.
- If we open up this box a little, we see that there are two parts to this mapping: analysis and synthesis.

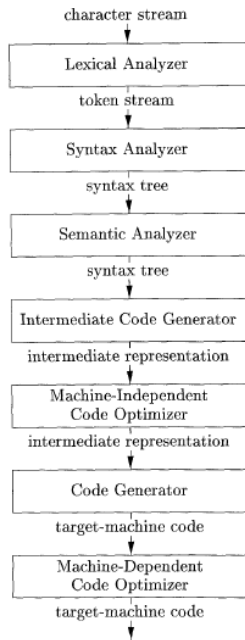
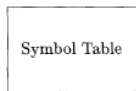


A compiler

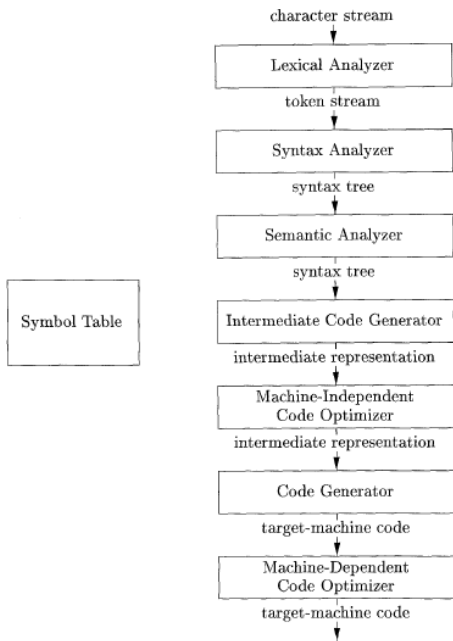
- The analysis part breaks up the source program into constituent pieces and imposes a grammatical structure on them.
- It then uses this structure to create an intermediate representation of the source program.
- If the analysis part detects that the source program is either syntactically ill formed or semantically unsound, then it must provide informative messages, so the user can take corrective action.



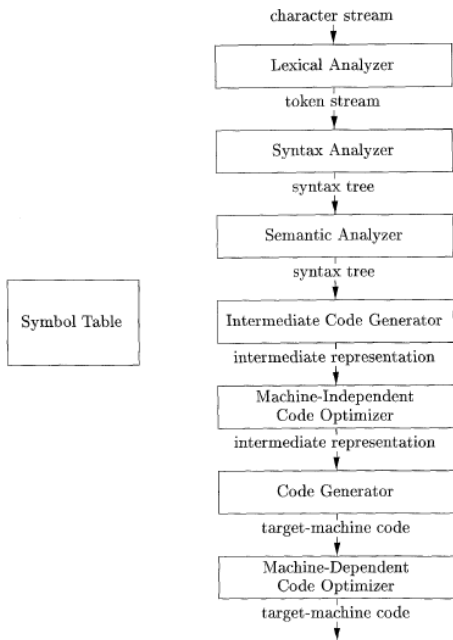
- The analysis part also collects information about the source program and stores it in a data structure called a symbol table, which is passed along with the intermediate representation to the synthesis part.
- The synthesis part constructs the desired target program from the intermediate representation and the information in the symbol table.
- The analysis part is often called the front end of the compiler; the synthesis part is the back end.



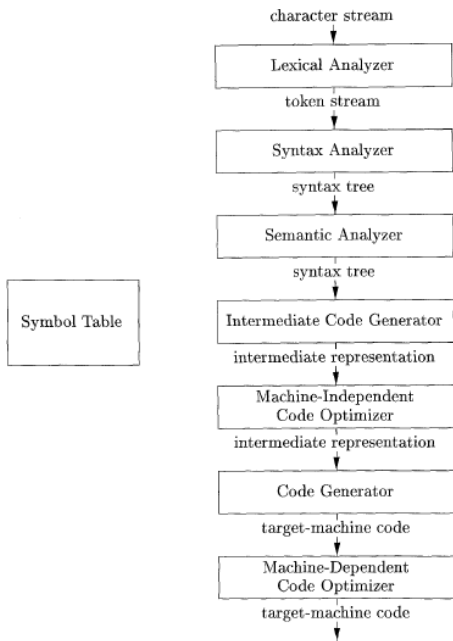
- If we examine the compilation process in more detail, we see that it operates as a sequence of phases, each of which transforms one representation of the source program to another.
- In practice, several phases may be grouped together, and the intermediate representations between the grouped phases need not be constructed explicitly.



- The symbol table, which stores information about the entire source program, is used by all phases of the compiler.

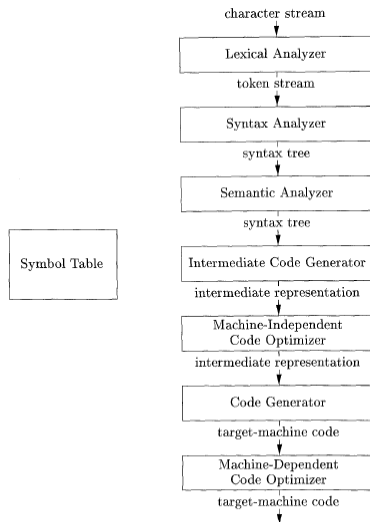


- Some compilers have a machine-independent optimization phase between the front end and the back end.
- The purpose of this optimization phase is to perform transformations on the intermediate representation, so that the back end can produce a better target program than it would have otherwise produced from an unoptimized intermediate representation.
- Since optimization is optional, one or the other of the two optimization phases shown may be missing.



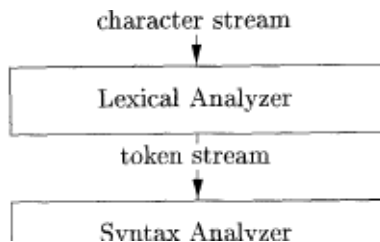
Lexical Analysis

- The first phase of a compiler is called lexical analysis or scanning.



Phases of a compiler

Lexical Analysis — *continued*



- The lexical analyzer reads the stream of characters making up the source program and groups the characters into meaningful sequences called **lexemes**.

Lexical Analysis — *continued*

- For each lexeme, the lexical analyzer produces as output a token of the form

<token-name, attribute-value>

that it passes on to the subsequent phase, syntax analysis.

- In the token, the first component **token-name** is an abstract symbol that is used during syntax analysis, and the second component **attribute-value** points to an entry in the symbol table for this token.
- Information from the symbol-table entry is needed for semantic analysis and code generation.

Lexical Analysis — *continued*

$\langle token-name, attribute-value \rangle$

`position = initial + rate * 60`

- We consider a source program that contains the assignment statement

`position = initial + rate * 60`

- The characters in this assignment could be grouped into the lexemes and mapped into the tokens passed on to the syntax analyzer

Lexical Analysis — *continued*

$\langle \text{token-name}, \text{attribute-value} \rangle$

`position = initial + rate * 60`

1. `position` is a lexeme that would be mapped into a token $\langle \text{id}, 1 \rangle$.
 - **id** is an abstract symbol standing for identifier,
 - 1 points to the symbol-table entry for `position`,
 - The symbol-table entry for an identifier holds information about the identifier, such as its name and type.

$\langle token-name, attribute-value \rangle$

`position = initial + rate * 60`

2. The assignment symbol = is a lexeme that is mapped into the token $\langle = \rangle$.
- Since this token needs no attribute-value, we have omitted the second component.
 - We could have used any abstract symbol such as assign for the token-name, but for notational convenience we have chosen to use the lexeme itself as the name of the abstract symbol.

Lexical Analysis — *continued*

$\langle \text{token-name}, \text{attribute-value} \rangle$

`position = initial + rate * 60`

3. `initial` is a lexeme that is mapped into the token $\langle \text{id}, 2 \rangle$,
 - 2 points to the symbol-table entry for `initial`.

Lexical Analysis — *continued*

$\langle token-name, attribute-value \rangle$

`position = initial + rate * 60`

4. `+` is a lexeme that is mapped into the token $\langle + \rangle$.

Lexical Analysis — *continued*

$\langle \text{token-name}, \text{attribute-value} \rangle$

`position = initial + rate * 60`

5. `rate` is a lexeme that is mapped into the token $\langle \text{id}, 3 \rangle$, where 3 points to the symbol-table entry for `rate`.

Lexical Analysis — *continued*

$\langle token-name, attribute-value \rangle$

`position = initial + rate * 60`

6. `*` is a lexeme that is mapped into the token $\langle * \rangle$.

Lexical Analysis — *continued*

$\langle token-name, attribute-value \rangle$

`position = initial + rate * 60`

7. 60 is a lexeme that is mapped into the token $\langle 60 \rangle$.

Lexical Analysis — *continued*

$\langle token-name, attribute-value \rangle$

`position = initial + rate * 60`

- Blanks separating the lexemes would be discarded by the lexical analyzer.

Lexical Analysis — *continued*

`position = initial + rate * 60`

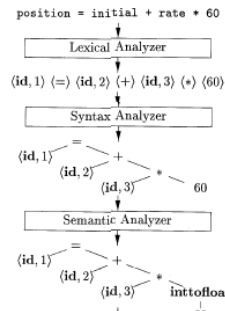
- Figure shows the representation of the assignment statement after lexical analysis as the sequence of tokens,

`<id, 1> <=> <id, 2> <+> <id, 3> <*> <60>`

- In this representation, the token names =, +, and * are abstract symbols for the assignment, addition, and multiplication operators, respectively.

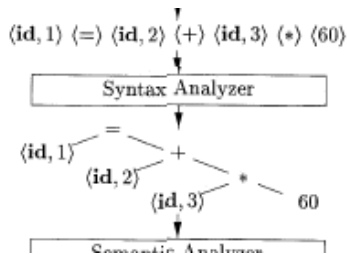
1	position	...
2	initial	...
3	rate	...

SYMBOL TABLE



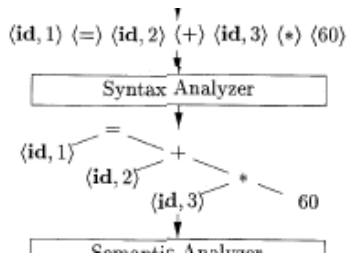
Syntax Analysis

- The second phase of the compiler is syntax analysis or parsing.
- The parser uses the first components of the tokens produced by the lexical analyzer to create a tree-like intermediate representation that depicts the grammatical structure of the token stream.



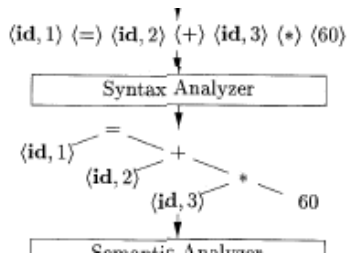
Syntax Analysis

- A typical representation is a syntax tree in which each interior node represents an operation and the children of the node represent the arguments of the operation.



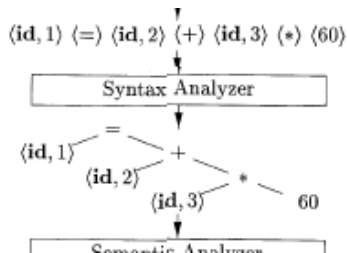
Syntax Analysis

- This tree shows the order in which the operations in the assignment
`position = initial + rate * 60`
are to be performed.



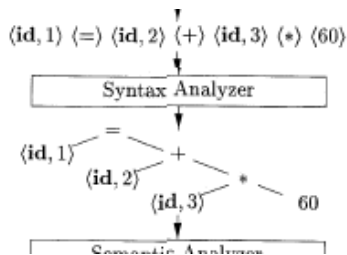
Syntax Analysis

- The tree has an interior node labeled $*$ with $\langle \text{id}, 3 \rangle$ as its left child and the integer 60 as its right child. The node $\langle \text{id}, 3 \rangle$ represents the identifier `rate`.
- The node labeled $*$ makes it explicit that we must first multiply the value of `rate` by 60.



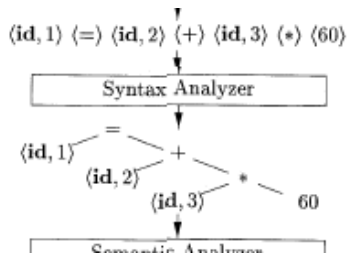
Syntax Analysis

- The node labeled `+` indicates that we must add the result of this multiplication to the value of `initial`.



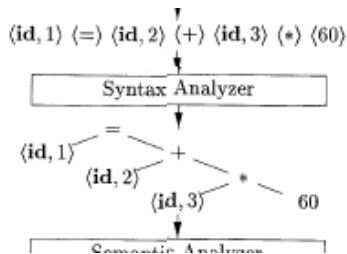
Syntax Analysis

- The root of the tree, labeled =, indicates that we must store the result of this addition into the location for the identifier `position`.



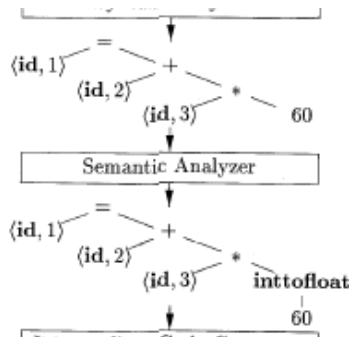
Syntax Analysis

- This ordering of operations is consistent with the usual conventions of arithmetic which tell us that multiplication has higher precedence than addition, and hence that the multiplication is to be performed before the addition.



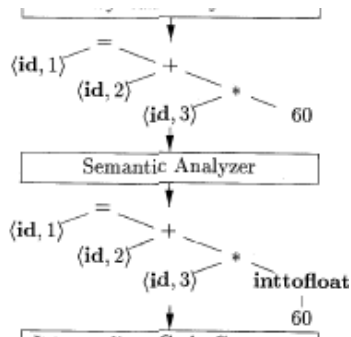
Semantic Analysis

- The semantic analyzer uses the syntax tree and the information in the symbol table to check the source program for semantic consistency with the language definition.
- It also gathers type information and saves it in either the syntax tree or the symbol table, for subsequent use during intermediate-code generation.



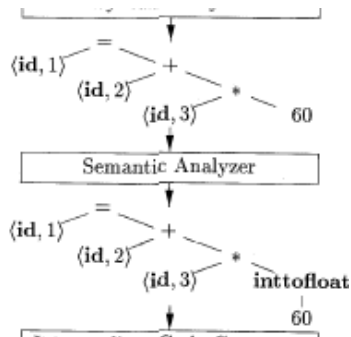
Semantic Analysis

- An important part of semantic analysis is type checking, where the compiler checks that each operator has matching operands.
- For example, many programming language definitions require an array index to be an integer.
- The compiler must report an error if a floating-point number is used to index an array.



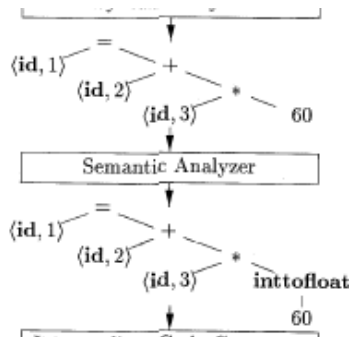
Semantic Analysis

- The language specification may permit some type conversions called coercions.
- For example, a binary arithmetic operator may be applied to either a pair of integers or to a pair of floating-point numbers.
- If the operator is applied to a floating-point number and an integer, the compiler may convert or coerce the integer into a floating-point number.



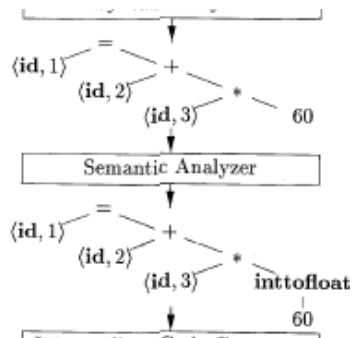
Semantic Analysis

- Suppose that `position`, `initial`, and `rate` have been declared to be floating-point numbers, and that the lexeme `60` by itself forms an integer.
- The type checker in the semantic analyzer discovers that the operator `*` is applied to a floating-point number `rate` and an integer `60`.
- In this case, the integer may be converted into a floating-point number.



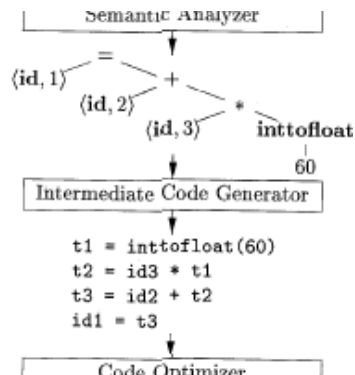
Semantic Analysis

- Notice that the output of the semantic analyzer has an extra node for the operator `inttofloat`, which explicitly converts its integer argument into a floating-point number.



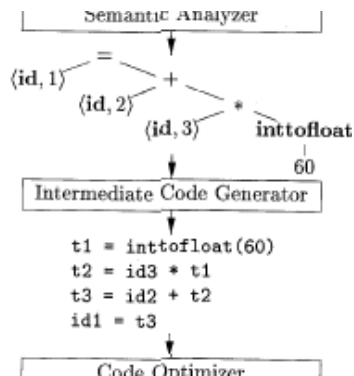
Intermediate Code Generation

- In the process of translating a source program into target code, a compiler may construct one or more intermediate representations, which can have a variety of forms.
- Syntax trees are a form of intermediate representation.
- They are commonly used during syntax and semantic analysis.



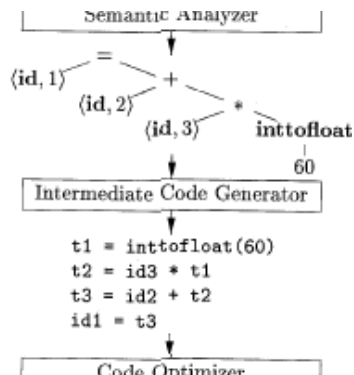
Intermediate Code Generation

- After syntax and semantic analysis of the source program, many compilers generate an explicit low-level or machine-like intermediate representation, which we can think of as a program for an abstract machine.
- This intermediate representation should have two important properties:
 - it should be easy to produce and
 - it should be easy to translate into the target machine.



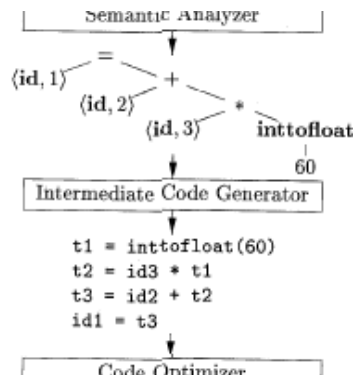
Intermediate Code Generation

- We consider an intermediate form called “three-address code”.
- Like the assembly language for a machine in which every memory location can act like a register.
- Three-address code consists of a sequence of instructions, each of which has at most three operands.



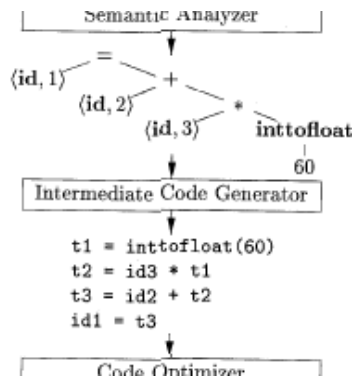
Intermediate Code Generation

```
t1 = inttofloat(60)
t2 = id3 * t1
t3 = id2 + t2
id1 = t3
```



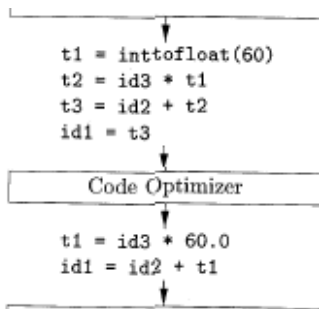
Intermediate Code Generation

- Each three-address instruction has at most one operator in addition to the assignment.
 - Has to decide on the order in which operations are to be done.
 - Multiplication precedes the addition in the source program.
- Must generate a temporary name to hold the value computed by each instruction.
- Some “three-address” instructions have fewer than three operands.



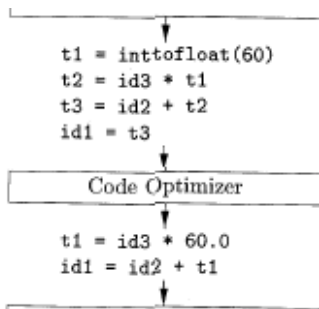
Code Optimization

- The machine-independent code-optimization phase attempts to improve the intermediate code so that better target code will result.
- Usually better means faster, but other objectives may be desired, such as shorter code, or target code that consumes less power.



Code Optimization

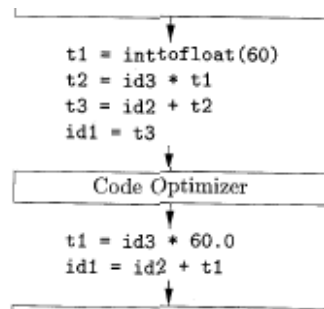
- A simple intermediate code generation algorithm followed by code optimization is a reasonable way to generate good target code.
- The optimizer can deduce that the conversion of 60 from integer to floating point can be done once and for all at compile time.
- So the `inttfloat` operation can be eliminated by replacing the integer 60 by the floating-point number 60.0.
- Moreover, `t3` is used only once to transmit its value to `<id, 1>`.



Code Optimization

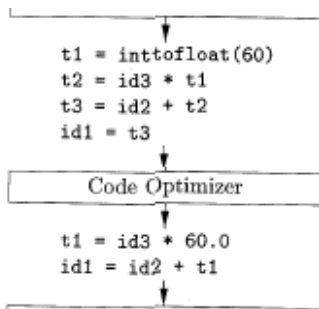
- So the optimizer can transform the code into the shorter sequence,

```
t1 = id3 * 60.0  
id1 = id2 + t1
```



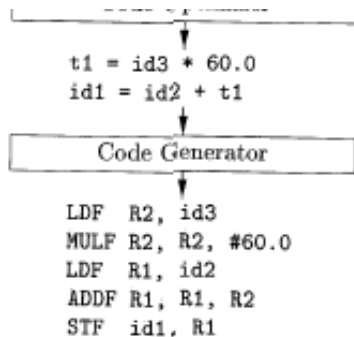
Code Optimization

- There is a great variation in the amount of code optimization different compilers perform.
- In those that do the most, the so-called “optimizing compilers,” a significant amount of time is spent on this phase.
- There are simple optimizations that significantly improve the running time of the target program without slowing down compilation too much.



Code Generation

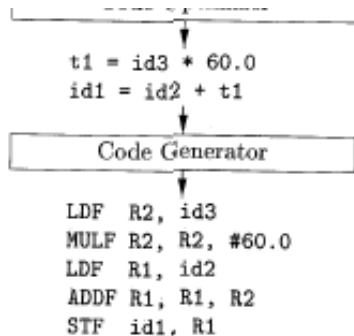
- The code generator takes as input an intermediate representation of the source program and maps it into the target language.
- If the target language is machine code, registers or memory locations are selected for each of the variables used by the program.
- Then, the intermediate instructions are translated into sequences of machine instructions that perform the same task.
- A crucial aspect of code generation is the judicious assignment of registers to hold variables.



Code Generation

- For example, using registers R1 and R2, the intermediate code might get translated into the machine code

```
LDF R2, id3
MULF R2, R2, #60.0
LDF R1, id2
ADDF R1, R1, R2
STF id1, R1
```



Symbol-Table Management

- An essential function of a compiler is to record the variable names used in the source program and collect information about various attributes of each name.

1	position	...
2	initial	...
3	rate	...

SYMBOL TABLE

Symbol-Table Management

- These attributes may provide information about the storage allocated for a name, its type, its scope (where in the program its value may be used)
- In the case of procedure names, such things as the number and types of its arguments, the method of passing each argument (for example, by value or by reference), and the type returned.

1	position	...
2	initial	...
3	rate	...

SYMBOL TABLE

Symbol-Table Management

- The symbol table is a data structure containing a record for each variable name, with fields for the attributes of the name.
- The data structure should be designed to allow the compiler to find the record for each name quickly and to store or retrieve data from that record quickly.

1	position	...
2	initial	...
3	rate	...

SYMBOL TABLE

The Grouping of Phases into Passes

- The discussion of phases deals with the logical organization of a compiler.
- In an implementation, activities from several phases may be grouped together into a pass that reads an input file and writes an output file.
- For example, the front-end phases of lexical analysis, syntax analysis, semantic analysis, and intermediate code generation might be grouped together into one pass.
- Code optimization might be an optional pass.
- Then there could be a back-end pass consisting of code generation for a particular target machine.

The Grouping of Phases into Passes

- The discussion of phases deals with the logical organization of a compiler.
- In an implementation, activities from several phases may be grouped together into a pass that reads an input file and writes an output file.
- For example, the front-end phases of lexical analysis, syntax analysis, semantic analysis, and intermediate code generation might be grouped together into one pass.
- Code optimization might be an optional pass.
- Then there could be a back-end pass consisting of code generation for a particular target machine.

The Grouping of Phases into Passes — *continued*

- Some compiler collections have been created around carefully designed intermediate representations that allow the front end for a particular language to interface with the back end for a certain target machine.
- With these collections, we can produce compilers for different source languages for one target machine by combining different front ends with the back end for that target machine.
- Similarly, we can produce compilers for different target machines, by combining a front end with back ends for different target machines.

The Grouping of Phases into Passes — *continued*

- Some compiler collections have been created around carefully designed intermediate representations that allow the front end for a particular language to interface with the back end for a certain target machine.
- With these collections, we can produce compilers for different source languages for one target machine by combining different front ends with the back end for that target machine.
- Similarly, we can produce compilers for different target machines, by combining a front end with back ends for different target machines.

The Grouping of Phases into Passes — *continued*

- Some compiler collections have been created around carefully designed intermediate representations that allow the front end for a particular language to interface with the back end for a certain target machine.
- With these collections, we can produce compilers for different source languages for one target machine by combining different front ends with the back end for that target machine.
- Similarly, we can produce compilers for different target machines, by combining a front end with back ends for different target machines.

Compiler-Construction Tools

- The compiler writer, like any software developer, can profitably use modern software development environments containing tools such as language editors, debuggers, version managers, profilers, test harnesses, and so on.
- In addition to these general software-development tools, other more specialized tools have been created to help implement various phases of a compiler.

Compiler-Construction Tools — *continued*

- These tools use specialized languages for specifying and implementing specific components, and many use quite sophisticated algorithms.
- The most successful tools are those that hide the details of the generation algorithm and produce components that can be easily integrated into the remainder of the compiler.

Compiler-Construction Tools — *continued*

Some commonly used compiler-construction tools include

1. **Parser generators** that automatically produce syntax analyzers from a grammatical description of a programming language.
2. **Scanner generators** that produce lexical analyzers from a regular-expression description of the tokens of a language.
3. **Syntax-directed** translation engines that produce collections of routines for walking a parse tree and generating intermediate code.

Compiler-Construction Tools — *continued*

Some commonly used compiler-construction tools include

4. **Code-generator generators** that produce a code generator from a collection of rules for translating each operation of the intermediate language into the machine language for a target machine.
5. **Data-flow analysis engines** that facilitate the gathering of information about how values are transmitted from one part of a program to each other part. Data-flow analysis is a key part of code optimization.
6. **Compiler-construction toolkits** that provide an integrated set of routines for constructing various phases of a compiler.



End of Slides