

Intro to Compilers

CSE 413, Autumn 2007
11-07-2007

11/07/2007 1

Agenda

- What's a compiler?
- Compilers vs. Interpreters
- Phases of a compiler

11/07/2007 2

And the point is...

- Execute this!


```
int nPos = 0;
int k = 0;
while (k < length) {
  if (a[k] > 0) {
    nPos++;
  }
}
```
- How?

11/07/2007 3

Compiler

- Read and analyze entire program
- Translate to semantically equivalent program in another language
 - Presumably easier to execute or more efficient
 - Should "improve" the program in some fashion
- Offline process
 - Tradeoff: compile time overhead (preprocessing step) vs execution performance

11/07/2007 4

Compiler vs. Assembler

```

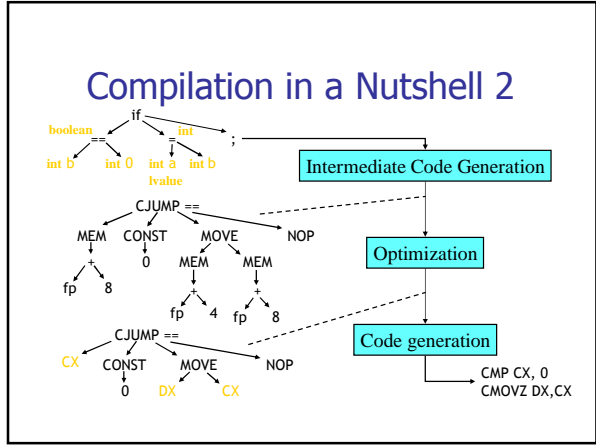
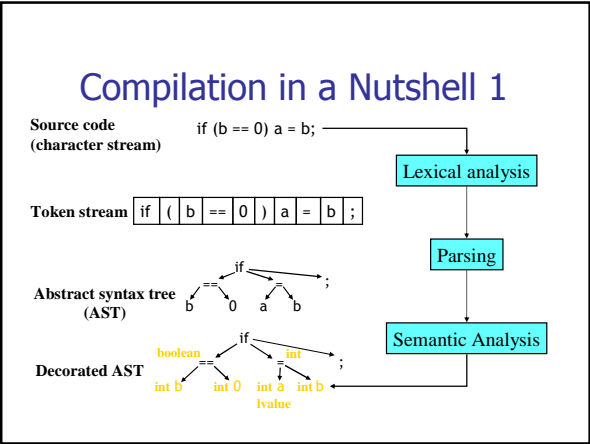
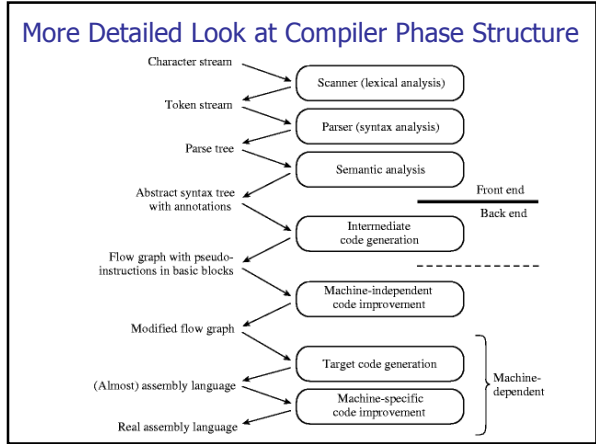
graph LR
  SP[Source program] --> C[Compiler]
  C --> AL[Assembly language]
  AL --> A[Assembler]
  A --> ML[Machine language]
  
```

11/07/2007 5

Assembler

- Principal tasks of an assembler are:
 - Replace opcodes and operands with their machine language encodings
 - Replace uses of symbolic names with actual addresses
- Assembler translates assembly language into Object code (Machine code)

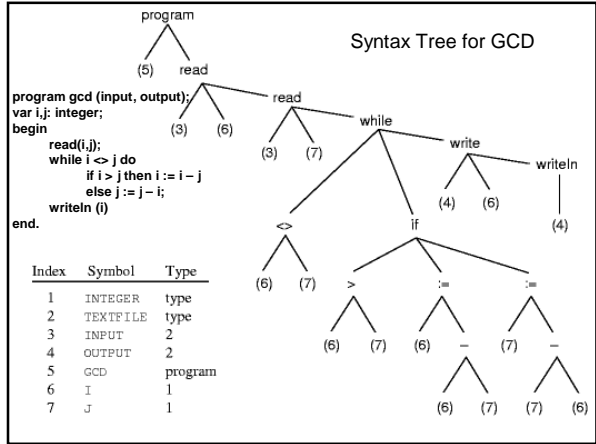
11/07/2007 6



```

program gcd (input, output);
var i,j: integer;
begin
  read(i,j);
  while i <> j do
    if i > j then i := i - j
    else j := j - i;
  writeln (i)
end.

```



- ### Compilers vs. Interpreters
- Interpreter**
 - A program that reads a source program and produces the results of executing that program
 - Compiler**
 - A program that translates a program from one language (the *source*) to another (the *target*)
- 11/07/2007 12

Common Issues

- Compilers and interpreters both must read the input – a stream of characters – and “understand” it; *analysis*

```
while(k < length){<nl> <tab> if(a[k] > 0
) <nl> <tab> <tab>{ n P o s + + ; } <nl> <tab> }
```

11/07/2007

13

Interpreter

- Interpreter
 - Execution engine
 - Program execution interleaved with analysis

```
running = true;
while (running) {
  analyze next statement;
  execute that statement;
}
```
 - May involve repeated analysis of some statements (loops, functions)

11/07/2007

14

Compiler

- Read and analyze entire program
- Translate to semantically equivalent program in another language
 - Presumably easier to execute or more efficient
 - Should “improve” the program in some fashion
- Offline process
 - Tradeoff: compile time overhead (preprocessing step) vs execution performance

11/07/2007

15

Typical Implementations

- Compilers
 - FORTRAN, C, C++, Java, COBOL, etc.
 - Strong need for optimization in many cases
- Interpreters
 - PERL, Python, Ruby, awk, sed, sh, csh, postscript printer, Scheme, Java VM
 - Effective if interpreter overhead is low relative to execution cost of individual statements

11/07/2007

16

Hybrid approaches

- Well-known example: Java
 - Compile Java source to byte codes – Java Virtual Machine language (.class files)
 - Execution
 - Interpret byte codes directly, or
 - Compile some or all byte codes to native code
 - Just-In-Time compiler (JIT) – detect hot spots & compile on the fly to native code
- Variation: .NET
 - Compilers generate MSIL
 - All IL compiled to native code before execution

11/07/2007

17

Why Study Compilers?

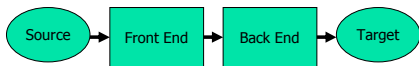
- Better Understanding Of Implementation Issues in Programming Languages:
 - How Is “This” Implemented?
 - Why Does “This” Run So Slowly?
- Translation appears several places:
 - Processing command line parameters
 - Converting files/programs from one language/format to another

11/07/2007

18

Structure of a Compiler

- First approximation
 - Front end: analysis
 - Read source program and understand its structure and meaning
 - Back end: synthesis
 - Generate equivalent target language program

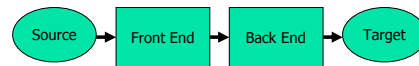


11/07/2007

19

Implications

- Must recognize legal programs (& complain about illegal ones)
- Must generate correct code
- Must manage storage of all variables
- Must agree with OS & linker on target format

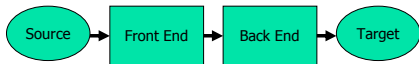


11/07/2007

20

More Implications

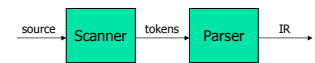
- Need some sort of Intermediate Representation(s) (IR)
- Front end maps source into IR
- Back end maps IR to target machine code
- May be multiple IRs – higher level at first, lower level in later phases



11/07/2007

21

Front End



- Split into two parts
 - Scanner: Responsible for converting character stream to token stream
 - Also strips out white space, comments
 - Parser: Reads token stream; generates IR
- Both of these can be generated automatically
 - Source language specified by a formal grammar
 - Tools read the grammar and generate scanner & parser (either table-driven or hard-coded)

11/07/2007

22

Tokens

- Token stream: Each significant lexical chunk of the program is represented by a token
 - Operators & Punctuation: {}[]!+ -=*;; ...
 - Keywords: if while return goto
 - Identifiers: id & actual name
 - Constants: kind & value; int, floating-point character, string, ...

11/07/2007

23

Scanner Example

- Input text


```
// this statement does very little
if (x >= y) y = 42;
```
- Token Stream

IF	LPAREN	ID(x)	GEQ	ID(y)
RPAREN	ID(y)	BECOMES	INT(42)	SCOLON

 - Notes: tokens are atomic items, not character strings; comments are *not* tokens

11/07/2007

24

Example

- Possible syntax for numeric constants

$digit ::= [0-9]$

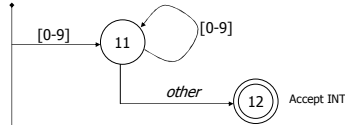
$digits ::= digit^+$

$number ::= digits (. digits) ?$
 $([eE] (+ | -) ? digits) ?$

11/07/2007

25

Scanner DFA Example



11/07/2007

26

Parser Output (IR)

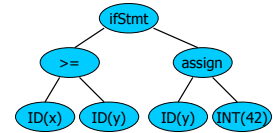
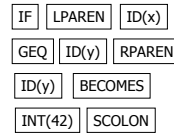
- Many different forms
 - Engineering tradeoffs that have changed over time
- Common output from a parser is an abstract syntax tree
 - Essential meaning of the program without the syntactic noise

11/07/2007

27

Parser Example

- Token Stream Input
- Abstract Syntax Tree



11/07/2007

28

Context-Free Grammars

- Formally, a grammar G is a tuple $\langle N, \Sigma, P, S \rangle$ where
 - N a finite set of non-terminal symbols
 - Σ a finite set of terminal symbols
 - P a finite set of productions
 - A subset of $N \times (N \cup \Sigma)^*$
 - S the *start symbol*, a distinguished element of N
 - If not specified otherwise, this is usually assumed to be the non-terminal on the left of the first production

11/07/2007

29

Grammar for a Tiny Language

- $program ::= statement \mid program \ statement$
- $statement ::= assignStmt \mid ifStmt$
- $assignStmt ::= id = expr ;$
- $ifStmt ::= if (expr) stmt$
- $expr ::= id \mid int \mid expr + expr$
- $Id ::= a \mid b \mid c \mid i \mid j \mid k \mid n \mid x \mid y \mid z$
- $int ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

11/07/2007

30

Example

```

program ::= statement | program statement
statement ::= assignStmt | ifStmt
assignStmt ::= id = expr;
ifStmt ::= if ( expr ) stmt
expr ::= id | int | expr + expr
id ::= a | b | c | i | j | k | n | x | y | z
int ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
  
```

G

w → a = 1 ; if (a + 1) b = 2 ;

11/07/2007 31

Static Semantic Analysis

- During or (more common) after parsing
 - Type checking
 - Check for language requirements like proper declarations, type compatibility
 - Preliminary resource allocation
 - Collect other information needed by back end analysis and code generation

11/07/2007 32

Back End

- Responsibilities
 - Translate IR into target machine code
 - Should produce fast, compact code
 - Should use machine resources effectively
 - Registers
 - Instructions
 - Memory hierarchy

11/07/2007 33

Back End Structure

- Typically split into two major parts with sub phases
 - "Optimization" – code improvements
 - May well translate parser IR into other IRs
 - We probably won't have time to do much with this part of the compiler, alas
 - Code generation
 - Instruction selection & scheduling
 - Register allocation

11/07/2007 34

The Result

- Input


```

if (x >= y)
  y = 42;
      
```
- Output


```

mov  eax,[ebp+16]
cmp  eax,[ebp-8]
jl   L17
mov  [ebp-8],42
L17:
      
```

11/07/2007 35