

CSc 553

Principles of Compilation

1 : Compiler Overview

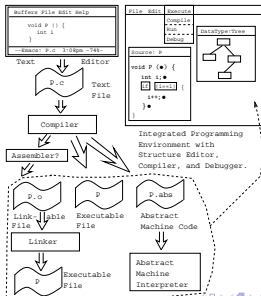
Department of Computer Science
University of Arizona

collberg@gmail.com

Copyright © 2011 Christian Collberg

What does a compiler do?

What's a Compiler???



Compiler Input and Output

Compiler Input

Text File Common on Unix.

Syntax Tree A structure editor uses its that knowledge of the source language syntax to help the user edit & run the program. It can send a syntax tree to the compiler, relieving it of lexing & parsing.

Compiler Output

Assembly Code Unix compilers do this. Slow, but easy for the compiler.

Object Code .o-files on Unix. Faster, since we don't have to call the assembler.

Executable Code Called a *load-and-go*-compiler.

Abstract Machine Code Serves as input to an *interpreter*. Fast turnaround time.

C-code Good for portability.

Static Semantic Analysis Is the program (statically) correct? If not, produce error messages to the user.

Code Generation The compiler must produce code that can be executed.

Symbolic Debug Information The compiler should produce a description of the source program needed by symbolic debuggers. Try `man gdb`.

Cross References The compiler may produce **cross-referencing** information. Where are identifiers declared & referenced?

Profiler Information The compiler should produce **profiler** information. Where does `my program` spend most of its execution time? Try `man gprof`.

The structure of a compiler

ANALYSIS

Lexical Analysis

Syntactic Analysis

Semantic Analysis

SYNTHESIS

Intermediate Code
Generation

Code Optimization

Machine Code
Generation

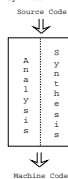
Compiler Organization

One Pass Analysis and Synthesis Fast. OK for definition-before-use languages like Pascal. No explicit intermediate representation. Target machine code is generated on-the-fly. Very little optimization is possible since we can't "look forward". Difficult to retarget, since semantic analysis and code generation are performed simultaneously.

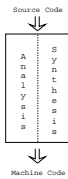
One Pass Plus Peephole Optimization Better code generation by performing a scan over the machine code and making local improvements.

One Pass Analysis + IR Generation Machine code is produced from an explicit intermediate representation. Better chances that the front-end & back-end can be recycled.

One Pass Analysis and Synthesis



One Pass plus Peephole Opt.



One Pass Anal. & IR Synth. + Code gen.

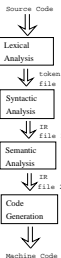


Multipass w/ Interm. Files Early compilers were severely constrained by the size of available primary storage. Therefore the compiler was often organized as a series of passes, where each pass wrote its output to an intermediate file which then became input to the next pass. Still a good design if you're not worried about speed.

Multipass Analysis Languages that allow "use-before-declaration", require the compiler to process the program more than once..

Multipass Synthesis Highly optimizing compilers usually process the intermediate representation in several passes. Often, we separate machine-independent and machine-dependent optimizations.

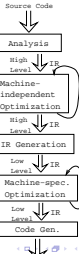
Multipass with multiple files

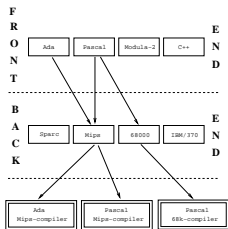


Multipass Analysis for fwd. ref.



Multipass Synthesis





Multipass Compilation

Multi-pass Compilation I

Multi-pass Compilation. . .

- We are going to work with compilers with multi-pass analysis and multi-pass synthesis parts.
- These compilers are very general:
 - They can handle any language, whether free or fixed declaration order.
 - They can produce efficient code.
 - They are portable since the front- and back-ends can be reused for compilers for new languages or new architectures.
- We will assume that the parser builds a tree (an **abstract syntax tree**) that is modified during semantic analysis, and then used during code generation.

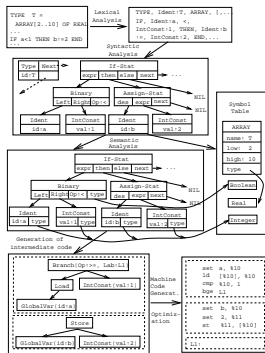
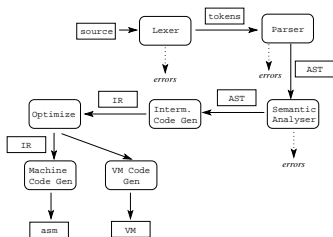
- The next slide shows the outline of a typical compiler. In a unix environment each pass could be a stand-alone program, and the passes could be connected by pipes:

```
lex x.c | parse | sem | ir | opt | codegen > x.s
```

- For performance reasons the passes are usually integrated:

```
front x.c > x.ir
back x.ir > x.s
```

The front-end does all analysis and IR generation. The back-end optimizes and generates code.



Example 1

Example

- Let's go through the compilation of a procedure Foo, from start to finish:

```

PROCEDURE Foo ();
VAR i : INTEGER;
BEGIN
  i := 1;
  WHILE i < 20 DO
    PRINT i * 2;
    i := i * 2 + 1;
  ENDDO;
END Foo;
  
```

- The compilation phases are:

Lexical Analysis ⇒ *Syntactic Analysis* ⇒
Semantic Analysis ⇒ *Intermediate code generation*
 ⇒ *Code Optimization* ⇒ *Machine code generation.*

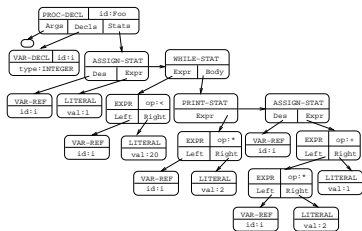
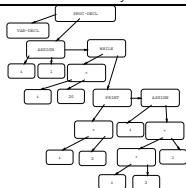
- Break up the source code (a text file) and into tokens.

Source Code	Stream of Tokens
PROCEDURE Foo ();	PROCEDURE, <id,Foo>, LPAR, RPAR, SC,
VAR i : INTEGER;	VAR, <id,i>, COLON, <id,INTEGER>,SC,
BEGIN	BEGIN, <id,i>,CEQ,<int,1>,SC,
i := 1;	WHILE, <id,i>, LT, <int,20>,DO,
WHILE i < 20 DO	PRINT, <id,i>, MUL, <int,2>, SC,
PRINT i * 2;	<id,i>, CEQ, <id,i>, MUL, <int,2>, PL
i := i * 2 + 1;	<int,1>, SC, ENDDO, SC, END, <id,Foo>
ENDDO;	
END Foo;	

Stream of Tokens

PROCEDURE, <id,Foo>, LPAR,RPAR,SC,VAR,<id,i>, COLON,<id,INTEGER>,SC, BEGIN,<id,i>,CEQ,<int,1>, SC,WHILE,<id,i>,LT,<int,20>, DO,PRINT,<id,i>,MUL,<int,2>, SC,<id,i>,CEQ,<id,i>,MUL, <int,2>,PLUS,<int,1>,SC, ENDDO,SC,END,<id,Foo>,SC

Abstract Syntax Tree

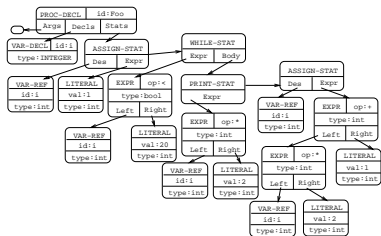


Abstract Syntax Tree

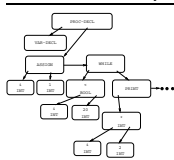


Decorated Abstract Syntax Tree





Decorated Abstract Syntax Tree



Intermediate Code

[1]	ASSIGN	i	1	
[2]	BRGE	i	20	[9]
[3]	MUL	t ₁	i	2
[4]	PRINT	t ₁		
[5]	MUL	t ₂	i	2
[6]	ADD	t ₃	t ₂	1
[7]	ASSIGN	i	t ₃	
[8]	JUMP	[2]		
[9]				

Example V/B – Intermediate Code Generation

Example VI – Code Optimization

Intermediate Code	Intermediate Code Definition
[1] ASSIGN i 1	ASSIGN A, B A := B;
[2] BRGE i 20 [9]	BRGE A, B, C IF (A ≥ B) THEN continue at instruction C;
[3] MUL t ₁ i 2	MUL A, B, C A := B * C;
[4] PRINT t ₁	ADD A, B, C A := B + C;
[5] MUL t ₂ i 2	SHL A, B, C A := shift B left C steps;
[6] ADD t ₃ t ₂ 1	PRINT A Print A and a newline;
[7] ASSIGN i t ₃	JUMP A Continue at instruction A;
[8] JUMP [2]	
[9]	

Intermediate Code

[1]	ASSIGN	i	1	
[2]	BRGE	i	20	[9]
[3]	MUL	t ₁	i	2
[4]	PRINT	t ₁		
[5]	MUL	t ₂	i	2
[6]	ADD	t ₃	t ₂	1
[7]	ASSIGN	i	t ₃	
[8]	JUMP	[2]		
[9]				

Optimized Intermediate Code

[1]	ASSIGN	i	1	
[2]	BRGE	i	20	[8]
[3]	SHL	t ₁	i	1
[4]	PRINT	t ₁		
[5]	ADD	t ₂	t ₁	1
[6]	ASSIGN	i	t ₂	
[7]	JUMP	[2]		
[8]				

Intermediate Code	MIPS Machine Code
	.data
	_i: .word 0
	.text
	.globl main
[1] ASSIGN i 1	main: li \$14, 1
[2] BRGE i 20 [8]	\$32: bge \$14, 20, \$32
[3] SHL t ₁ i 1	sll \$a0, \$14, 1
[4] PRINT t ₁	li \$v0, 1
[5] ADD t ₂ t ₁ 1	syscall
[6] ASSIGN i t ₂	addu \$14, \$a0, 1
[7] JUMP [2]	b \$32
[8]	\$33: sw \$14, _i

Summary

Readings and References

Summary I

- Read the Dragon Book:
 - Introduction Chapter 1
 - A Simple Syntax-Directed Translator Chapter 2
 - A Complete Front-End Appendix A

- The structure of a compiler depends on
 - 1 the complexity of the language we're working on (higher complexity \Rightarrow more passes),
 - 2 the quality of the code we hope to produce (better code \Rightarrow more passes),
 - 3 the degree of portability we hope to achieve (more portable \Rightarrow better separation between front- and back-ends).
 - 4 the number of people working on the compiler (more people \Rightarrow more independent modules).
- Some highly retargetable compilers for high-level languages produce C-code, rather than machine code. This C-code is then compiled by the native C compiler to machine code.

- Some languages (APL, LISP, Smalltalk, Java, ICON, Perl, Awk) are traditionally *interpreted* (executed in software by an *interpreter*) rather than compiled to machine code.
- Some interpreters use *dynamic compilation (or jitting)*, switching between
 - 1 interpreting the virtual machine code,
 - 2 translating the virtual machine code to native machine code,
 - 3 executing the native machine code,
 - 4 optimizing the native and/or virtual machine code, and
 - 5 throwing native code away if it is no longer needed or takes up too much room.

All this is done dynamically at runtime.

Historical Notes

The First Compiler

- FORTRAN I was the first "high-level" programming language. It's designers also wrote the first real compiler and invented many of the techniques that we use today.
- The FORTRAN manual can be found here:
<http://www.fh-jena.de/~kleine/history>.
- The excerpt on the next few slides is taken from
John Backus, The history of FORTRAN I, II, and III, History of Programming Languages, The first ACM SIGPLAN conference on History of programming languages, 1978.

Before 1954 almost all programming was done in machine language or assembly language. Programmers rightly regarded their work as a complex, creative art that required human inventiveness to produce an efficient program. Much of their effort was devoted to overcoming the difficulties created by the computers of that era: the lack of index registers, the lack of builtin floating point operations, restricted instruction sets (which might have AND but not OR, for example), and primitive input- output arrangements. Given the nature of computers, the services which "automatic programming" performed for the programmer were concerned with overcoming the machine's shortcomings. Thus the primary concern of some "automatic programming" systems was to allow the use of symbolic addresses and decimal numbers...

