

50006 - Compilers - Lecture 1

Oliver Killane

05/01/22

Course Introduction

Lecturers

- **Paul Kelly**
 - Introduction to syntax analysis & toy compiler
 - Code Generation
 - Generating better code using registers
 - Register Allocation
 - Optimisation and data flow analysis
 - Loop optimisations
- **Naranker Dulay**
 - Lexical Analysis (characters to tokens)
 - LR (Bottom-up) parsers (tokens to AST)
 - LL (top-down parsing)
 - Semantic Analysis, checking program validity
 - Runtime organisation (memory)

Materials

- Scientia
- Course Textbook
- Course Webpages for Paul Kelly and Naranker Dulay

Recommended Textbooks:

- **The Dragon Book: 'Compilers: Principles, Techniques and Tools'**
The definitive book on compilers. Called the dragon book due to its front illustration (allusion to taming the dragon of compiler complexity).

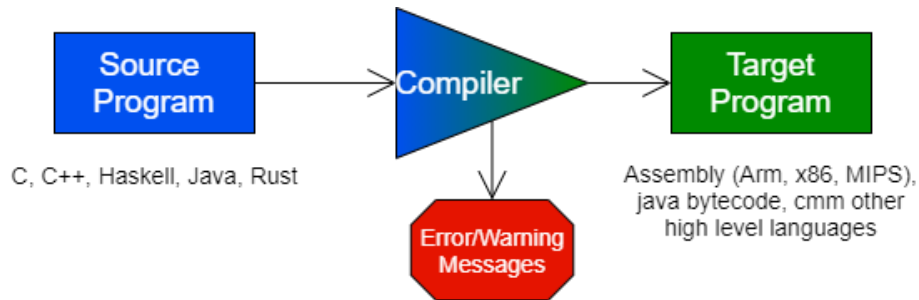
Has all required information and lots more (very thick book). ([Amazon link](#))
- **'Modern Compiler Implementation in Java'**
Details a java compiler implementation, with reasoning behind design choices. ([Amazon link](#))
- **Most Recommended: 'Engineering a Compiler'**
Most closely follows course & well reviewed. ([Amazon link](#))

Compilers Introduction

A particular class of program called **language processors**. Compilers can be:

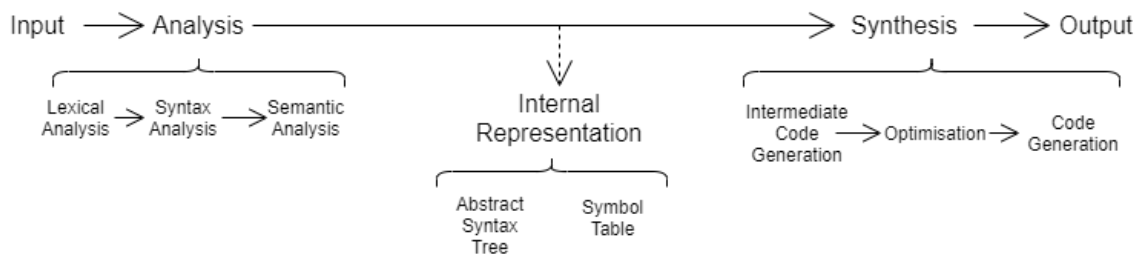
- **Processes** programs written in some language.
- **Writes** programs in some language.
- **Translates** programs from one language to an equivalent program in another. (e.g from high to lower level language)

A tool to enable programs written using high-level concepts to be translated into a low-level implementation.



- Typically from a high level language to a lower level one.
- Error messages ensure source program can be fixed easily when failing to compile.
- Can compiler from & to a high-level language (e.g compile to C, or between languages)
- Analyse the source program to provide the programmer useful information.

Basic Compiler Structure



- **Input** Take program written in some language.
- **Analysis**
Create an internal representation (IR) of the source which encodes its meaning (semantics).

This often starts with a tree (abstract syntax tree), but graphs may be used to represent control flow.
 - **Lexical Analysis**
Tokenization, converting text to a series of tokens representing keywords, identifiers, etc.
 - **Syntax Analysis**
Analysis of the nested structure of a program (e.g if-else within if-else within loop block).
 - **Semantic Analysis**
Analysing the meaning of a program such as type checking, determining overloading (e.g is + between variables correct & what action (e.g add integers, or concantenate strings)).
- **Synthesis**
Use the IR to create the program in the target language (retaining the same meaning/seman-tics).

The IR may be analysed and transformed to provide extra information to the programmer

(e.g warnings) as well as to optimise (e.g inlining, loop hoisting, loop unrolling, removing dead code).

- **Intermediate Code Generation**

Encodes more information, and is used in sophisticated compilers for optimisation. For example determinign loop invariant code (code that can be hoisted out of the loop & only be evaluated once).

- **Optimisation**

Improving code performance without altering its meaning.

- **Output** Output the target program. (e.g x86 assembly or Java bytecode)

Symbol Table

Contains information on all declared identifiers, and is used to check its use (e.g type checking) and generate code for access.

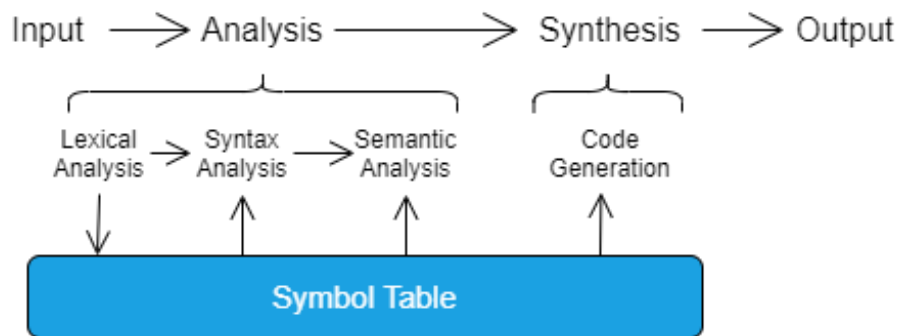
- Name
- Type of identifier (function, class, variable, constant)
- Size (Memory to reserve & where (e.g stack, data segment, heap))

Code	GCC symbol table
1 <code>#include <stdio.h></code>	1 00000000000003dc8 d _DYNAMIC
2 <code>#include <stdlib.h></code>	2 00000000000003fb8 d _GLOBAL_OFFSET_TABLE_
3	3 00000000000002000 R _IO_stdin_used
4	4 w _ITM_deregisterTMCloneTable
5	5 w _ITM_registerTMCloneTable
6 <code>int a = 6;</code>	6 000000000000021b4 r _FRAME_END_
7 <code>void pretty_print(int num) {</code>	7 00000000000002018 r _GNU_EH_FRAME_HDR
8 <code>printf("The number is: %d", num);</code>	8 00000000000004018 D _TMC_END_
9 <code>}</code>	9 00000000000004014 B __bss_start
10	10 w __cxa_finalize@@GLIBC_2.2.5
11 <code>int another_cool_function(int* num) {</code>	11 00000000000004000 D __data_start
12 <code>*num++;</code>	12 00000000000001100 t __do_global_ctors_aux
13 <code>return *num - 1;</code>	13 00000000000003dc0 d ↪ __do_global_ctors_aux_fini_array_entry
14 <code>}</code>	14 00000000000004008 D __dso_handle
15	15 00000000000003db8 d ↪ __frame_dummy_init_array_entry
16 <code>int main(int argc, char **argv) {</code>	16 w __gmon_start__
17 <code>int b = 7;</code>	17 00000000000003dc0 d __init_array_end
18 <code>pretty_print(a + b);</code>	18 00000000000003db8 d __init_array_start
19 <code>return EXIT_SUCCESS;</code>	19 00000000000001240 T __libc_csu_fini
20 <code>}</code>	20 000000000000011d0 T __libc_csu_init
21	21 U __libc_start_main@@GLIBC_2
22	22 ↪ .2.5
	23 00000000000004014 D _edata
	24 00000000000004018 B _end
	25 00000000000001248 T _fini
	26 00000000000001000 t _init
	27 00000000000001060 T _start
	28 00000000000004010 D a
	29 00000000000001171 T another_cool_function
	30 00000000000004014 b completed.8060
	31 00000000000004000 W data_start
	32 00000000000001090 t deregister_tm_clones
	33 00000000000001140 t frame_dummy
	34 00000000000001194 T main
	35 00000000000001149 T pretty_print
	36 U printf@@GLIBC_2.2.5
	37 00000000000001c0 t register_tm_clones

GCC Symbol Table

Simply compile, and then use `nm` on the resulting binary to get a basic representation of GCC's symbol table (that is included with the executable).

The symbol table is used in different ways through different phases of the compiler:



Code

```

1 int a;
2 int b;
3
4 void test_fun(void) {
5     a = b+42;
6 }

```

Boilerplate

Not lots of boilerplate assembly is omitted, if you try this locally you will see this.

Assembly

```

1 test_fun:
2     movl    b(%rip), %eax
3     addl    $42, %eax
4     movl    %eax, a(%rip)
5     ret
6 b:
7     .zero   4
8 a:
9     .zero   4

```

GodBolt

The compiler explorer is a great tool for generating & exploring assembly generated by compilers.

Compiler Phases

We can show the phases of the compiler through an example. Converting **C** code into **AT&T/GAS** assembly.

```
int_space a; newline int_space b; newline newline void_space test_fun(void)_space { newline a_space =_space b+42; newline }
```

```

[INTtok,
IDENTtok(pointer to 'a' entry in symbol table),
SEMICOLONtok,
INTtok,
IDENTtok(pointer to 'b' entry in symbol table),
SEMICOLONtok,
VOIDtok,
IDENTtok(pointer to 'test_fun' entry in symbol table),
LBRACKETtok,
VOIDtok,
RBRACKETtok,
LCURLYBRACKETtok,
IDENTtok(pointer to 'a' entry in symbol table),

```

EQtok,
IDENTtok(pointer to 'b' entry in symbol table),
PLUStok,
CONSTINTtok(123),
SEMICOLONtok,
RCURLYBRACKETtok]

Lexical Analysis

Also known as tokenization. The input sequence of characters is tokenized and any new identifiers added to the symbol table during tokenization.

Tokens to identifiers include pointers to their corresponding entry in the symbol table.

Syntax Analysis

Also known as parsing. Using a grammatical structure to generate an **Abstract Syntax Tree** from the sequence of tokens provided by the previous stage.

1. Ensure the structure is correct, report an error otherwise.
2. Build **Abstract Syntax Tree** from tokens to represent program.

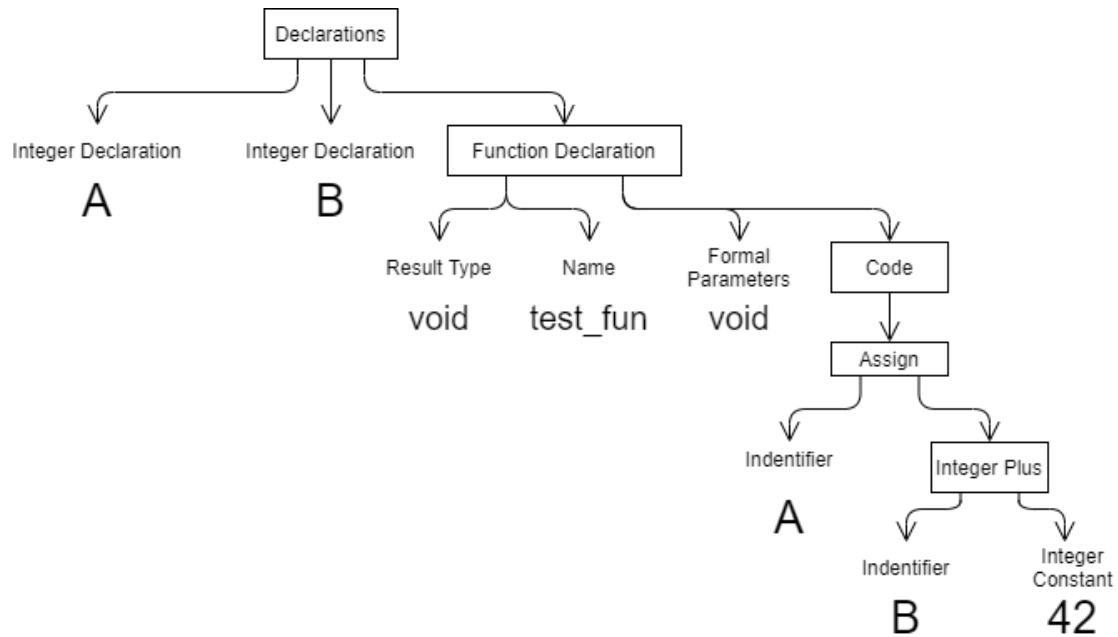
The grammar of the language can be expressed in a notation such as **Backus-Naur Form**. This is much the same as the definitions used for the simple while language in **Models of Computation**.

A rudimentary example could be:

$$\begin{aligned} Cond &::= 'true' \mid 'false' \mid Expr == Expr \mid Expr < Expr \mid Cond \& Cond \mid \neg Cond \dots \\ Expr &::= Var \mid Const \mid Expr + Expr \mid Expr \times Expr \dots \\ Stat &::= x := Expr \mid 'if' Cond 'then' Stat 'else' Stat \mid Stat; Stat \mid 'skip' \mid 'while' Cond 'do' Stat \end{aligned}$$

The **Abstract Syntax Tree** is implemented as a tree of linked objects. It must be designed to be efficient to construct & contain all require information for subsequent stages.

The **Abstract Syntax Tree** for the basic C program would be something like:



Domain Specific Languages

Many programming languages are **domain specific**, meaning they are designed for a specific use.

- **Latex** Creating documents
- **Markdown** Simple text files with image, font support (e.g readmes)
- **YAML & TOML** Configuration files
- **Verilog** Digital Circuit Design
- **R** Statistics
- **Simulink** Dynamic System Modelling
- **Yacc** (Yet another compiler-compiler) parser generator language

There are also many language processors that are not compilers:

- **FindBugs** Finds bugs in java code
- **Pylint** Linter for python
- **Valgrind** Detecting bugs & memory leaks