

Practical Compiler Development Tutorial – Part 1

This is a rather informal introduction to development of a hobby compiler . The more formal chapters on compiler development will be given in later tutorials. By the end of this tutorial you will be able to create a simple interpreter . This can be easily converted into a compiler. The requirements for following this tutorial are

1. A reasonable grasp of C / C++ / Python / C# / Java
2. An elementary knowledge of basic data structures and recursion

Compiler defined

A compiler is defined as a program that converts a source file in a specified form to an output file that can be fed to an assembler or simulated by a machine (virtual / real). Most of the compiler convert the given source file into assembly language that can be assembled by an assembler of the target machine.

For example to see the code generated by your source program in Turbo C++ , type the following

```
tcc -S source.c
```

Phases of a compiler

Lexical analysis : The lexical analysis phase of the compiler essentially does the tokenizing . This phase takes a stream of text as input and converts it into tokens required for the parser . For example the lexical analyzer splits the C source into if , else , while , int etc . The lexical analyzer can be constructed using a hard coded dfa or by using a generator like lex . In the next part of this tutorial I will describe how to write your own lex like generator . In this tutorial we will use an ad hoc lexical analyzer which will be sufficient for the purpose of demonstration. In fact you will be surprised to see that this sort of approach is taken in small hobby compilers like small- C .

Syntax Analysis : The lexemes must be arranged in such form that reflects the syntactic structure of the language . This is called syntax analysis or parsing . In this tutorial we will discuss only a form of top down parsing known as recursive descent parsing with one look ahead . Every programming language has a grammar which is usually represented in BNF form . A CFG grammar of a programming language is defined by the following

1. A set of terminal symbols that cannot be substituted
2. A set of non terminals which can be substituted by terminal or other non terminals
3. A set of Productions that define the grammar (of the form $NT \rightarrow A \mid B$, ie LHS should have only one nonterminal)
4. A start symbol from which everything begins

Eg A grammar for the string aaaa ..any number of times is

$S \rightarrow "a" S$ { terminals are enclosed in "" and non terminals in capital letters }
 { | used to denote OR }

$S \rightarrow \text{nothing}$

ie this means that S can be an "a" followed by S itself or S can be nothing .

Lets see how aaa is derived

$S \rightarrow aS \Rightarrow S \rightarrow aaS$ (using $S \rightarrow aS$) $\Rightarrow S \rightarrow aaaS$ (using $S \rightarrow aS$) $\Rightarrow S \rightarrow aaa$ (using $S \rightarrow \text{nothing_left}$) .

Now lets see how a recursive descent parser is constructed for this grammar. This is easily done by replacing each non terminal by function call and performing a match for the terminal

```
Void parseS()
{
    match_current_input_token_for("a");
    if( not_matched) error();
    if(end_of_input) return ;
    parseS();
}
```

In order to construct a recursive descent parser like the above the grammar should follow these conditions .

For each production of the form $S \rightarrow A \mid C \mid E$.

The first sets of each non-terminal in each of this production must be disjoint
ie $FIRST(A) \cap FIRST(C) \cap FIRST(E) = \text{NULL}$

The first sets and the follow sets of a nonterminal must be disjoint ie ...

$FIRST(A) \cap FOLLOW(A) = \text{NULL}$

$FIRST(C) \cap FOLLOW(C) = \text{NULL}$

...so on ...

This implies that the first set should not be equal to the follow set of the non- terminal.

The first set of a terminal is the symbol itself. The first set of the non – terminal is the set of terminals with which the nonterminal can begin with . The follow set of a non terminal is the set of terminals which can come after the non terminal. Lets find out the First n Follow for the following grammar .

$S \rightarrow "a" A \mid C D$
 $A \rightarrow "b" A \mid "c"$
 $C \rightarrow "e" \mid \text{nothing}$
 $D \rightarrow "d" \mid "g" D$

$FIRST(S) = FIRST("a" A) \cup FIRST(CD)$
 $= "a" \cup FIRST(C)$
 $= "a" \cup FIRST(C) \cup FIRST(D)$ { since C can be nothing , first of D must be included}

 $= "a" \cup FIRST("e") \cup FIRST("d") \cup FIRST("g")$
 $= "a" \cup "e" \cup "d" \cup "g"$
 $= \{ a , e , d , g \}$

Similarly $FOLLOW(C) = FIRST(D) = \{ d , g \}$ etc ... Computing the first and follow of the rest of the symbols in the grammar is left as an exercise to you . These rules have to followed in order to avoid conflicts in the predictive parsing table. More about advanced parsing methods in later chapters .

Now i present a more formal algorithms for finding first and follow .is given below . You might find it easier to follow the informal definitions that the algorithms given below.

Finding all nothingable non terminals

- 1) if given A -> nothing the set NOTINGABLE(A) = true
- 2) while no changes made in NOTINGABLE table
- 3) if production such that A -> A1 A2 A3 ... Ak and all of NOTHINGABLE(A1) NOTHINGABLE(A2) ... NOTHINGABLE(Ak) = true
- 4) then NOTINGABLE(A) = true

Finding first

- 1) if terminal then set FIRST (Symbol) = symbol
- 2) while no changes made in FIRST table
- 3) if production such that A -> A1 A2 A3 ... Ak
- 4) FIRST(A) = FIRST(A) union FIRST (A1)
- 5) while (NOTHINGABLE(A1) , NOTHINGABLE(A2) .. NOTHINGABLE(Aj). = TRUE)
- 6) FIRST(A) = FIRST(A) U FIRST(A1) U FIRST (A2) U FIRST (A3) ... FIRST(Aj+1)

Finding follow

- 1) follow is not defined for terminals
- 2) follow of the Start symbol is the end of input marker
- 3) if there exists a nonterminal A such that B -> E A D then
- 4) add FIRST(D) to FOLLOW(A)
- 5) if there exists a production such that A -> G H such that NOTHINGABLE(H) = true then
- 6) everything in FOLLOW (A) is also in FOLLOW (H)

As a final example, grammar for if statement.

S -> IFCONDTION

IFCONDITON -> “if” EXPRESSION “then” STATEMENT

EXPRESSION -> TERM “+” TERM

TERM -> FACTOR *FACTOR

FACTOR -> id

EXPRESSION -> (“ EXPRESSION “)”

STATEMENT -> id “,” STATEMENT | nothing

It can be seen that right recursion in many cases is unneeded and can be removed using a simple loop .
As right recursion leads to tail recursion. Eg

The grammar $S \rightarrow aS \mid \text{nothing}$ can be more compactly written as $S \rightarrow (a)^*$

The parser reduces to

```
parseS()
{
    while(current_token == "a" ) advance_input_pointer();
}
```

Left recursion will lead to an infinite loop and should be removed .For eg the grammar

$S \rightarrow SB \mid D$

can be re written as

$S \rightarrow D S_DASH$

$S_DASH \rightarrow B S_DASH$

Common factor in grammars should also be removed eg for the grammar

$S \rightarrow \text{"if" condition "then" statement}$

$S \rightarrow \text{"if" condition "then" statement "else" statement}$

can be modified into

$S \rightarrow AB$

$A \rightarrow \text{"if" condition "then" statement}$

$B \rightarrow \text{"else" statement} \mid \text{nothing}$

Semantic Analysis : Syntax analysis does not tell anything about the meaning of the statement under consideration . It simply check whether the syntax is correct . For eg

```
const int i = 20;
```

```
i++;
```

The above lines are syntactically correct but semantically wrong . The semantics are sometimes given using attribute grammars . .More about this later.

Code Generation : This is the final phase where the assembly code to the target machine is generated . I have skipped many other phases which are part of a more mature compiler. The purpose of this tutorial is to give quick introduction . The most of the hobby compilers the major component is the parser and code generation and semantic analysis is done when each non terminal is encountered in the parser. As an example of developing a simple interpreter , I have provided the source code of a simple basic like interpreter. I have used the mingw port of gcc for compiling the source . I used dev-cpp as the ide . The corresponding project and source files are also submitted .

