

Syntax Analysis: Scanning and Parsing

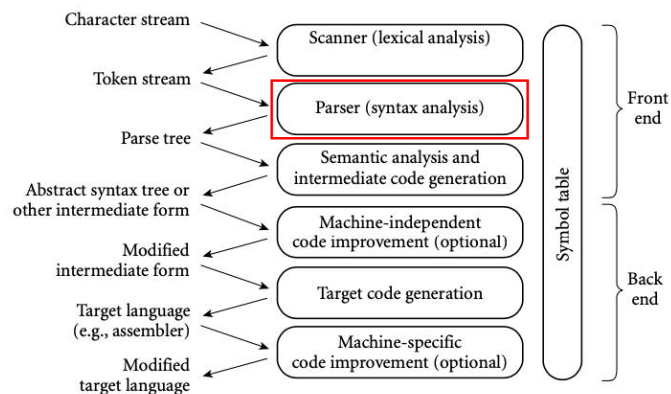
204315: OPL

Outline

- Scanner
 - Regular Expression
 - Deterministic Finite Automaton
- Parser
 - Context-Free Grammar
 - LL Parsing
 - LR Parsing

2

Compilation Process Overview (Recap)

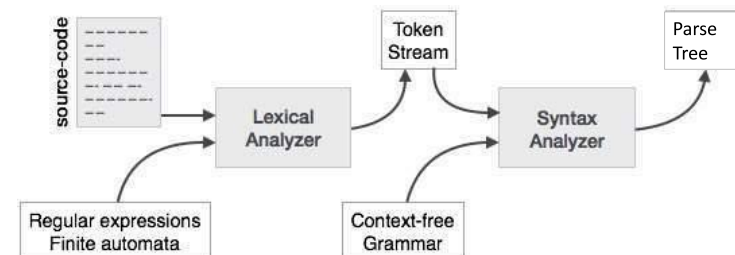


Scott, M. L. (2016). *Programming Language Pragmatics*.

3

Syntax Analyzer

- **Syntax analyzer or parser**
 - takes the input from a lexical analyzer in the form of token streams
 - analyzes the source code (token stream) against the rules to detect any **syntax errors** in the code.
 - The output of this phase is a parse tree.



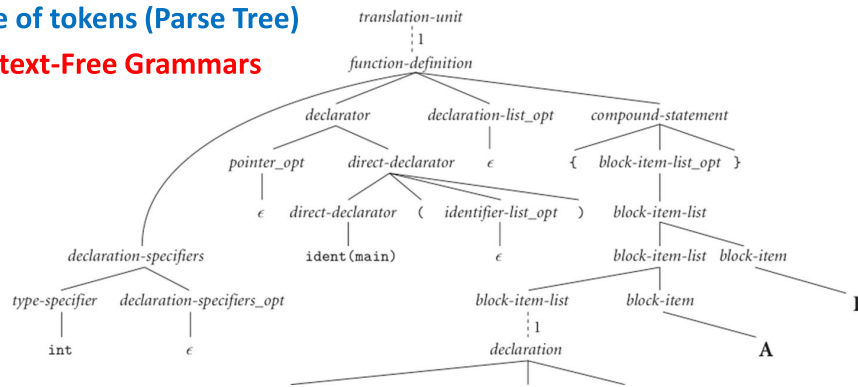
https://www.tutorialspoint.com/compiler_design/compiler_design_syntax_analysis.htm

4

Parsing phase

- Build tree of tokens (Parse Tree)

- Uses **Context-Free Grammars**



Scott, M. L. (2016). *Programming Language Pragmatics*.

5

Limitation of Regular Expression

- REs work well for defining tokens such as identifiers, constant etc..
- But unable to specify nested structure which are central to programming languages

- Therefore, cannot generate

- balanced parentheses
- nested chain structures



```
struct student_college_detail
{
    int college_id;
    char college_name[50];
};

struct student_detail
{
    int id;
    char name[20];
    float percentage;
    // structure within structure
    struct student_college_detail clg_data;
}stu_data;
```

[https://www.quora.com/What-is-a-nested-structure-and-explain-with-syntax](https://www.quora.com/What-is-a-nested-structure-and-explain-with-examples-with-syntax)

https://www.tutorialspoint.com/compiler_design/compiler_design_syntax_analysis.htm

6

Context-Free Grammars: Overview

- Context-free grammars (CFGs)
 - a set of recursive rules used to generate patterns of strings
 - used to describe context-free languages
 - CFG is sometimes called Backus-Naur Form (BNF)
- CFSs are studied in fields of
 - theoretical computer science
 - compiler design
 - linguistics
- CFG's are used to
 - describe programming languages and construct parser program in compiler

Karleigh Moore, Alex Chumbley, and Jimin Khim, Context Free Grammars: <https://brilliant.org/wiki/context-free-grammars/>

7

Context-Free Grammars: Formal definition

- A context-free grammar G is defined by the 4-tuple $G = \{V, \Sigma, R, S\}$
 - V is a set of non-terminals
 - represents a different type of phrase or clause in the sentence.
 - Σ , is a finite set of terminals, disjoint from V
 - represent actual content of the sentence
 - R , is a set of production rules
 - Member of R is a relation from $V \rightarrow V \times \Sigma$
 - S , is a start symbol S
 - Represents the whole program

8

Context-Free Grammars: Example

- $V = \{expr, op\}$ #non-terminal
 - $\Sigma = \{id, number, +, -, *, /, (,)\}$ #terminal
 - $R =$

$expr \rightarrow id \mid number \mid - expr \mid (expr)$
 $\quad \mid expr op expr$
 $op \rightarrow + \mid - \mid * \mid /$

#production rules
 - $S = expr$ # start symbol
- Each of the rules in a CFG is **known as a production**.
 - The **left-hand side symbols** of the productions are known as **variables**, or **non-terminals**, e.g., the language's tokens.
 - **Terminals cannot appear on the left-hand side of any productions.**
 - The **left-hand side** of the first production, is called the **start symbol**. It names the construct defined by the overall grammar.

9

CFG usage

$expr \rightarrow id \mid number \mid - expr \mid (expr)$
 $\quad \mid expr op expr$
 $op \rightarrow + \mid - \mid * \mid /$

- CFG can be used to check if a string is grammatical (with respect to the language)
- Can the following statement be derived from the grammars on the top ?

slope * x + intercept

10

Checking steps

Rule 1

$expr \rightarrow id \mid number \mid - expr \mid (expr)$
 $\quad \mid expr op expr$

Rule 2

$op \rightarrow + \mid - \mid * \mid /$

- Repeatedly rewrite the right-most-term until everything is terminal symbol
- Derivation steps for "slope * x + intercept" are

$expr \Rightarrow expr op expr$
 $\Rightarrow expr op id$
 $\Rightarrow expr + id$
 $\Rightarrow expr op expr + id$
 $\Rightarrow expr op id + id$
 $\Rightarrow expr * id + id$
 $\Rightarrow id * id + id$
 $\quad (slope) \quad (x) \quad (intercept)$

Rule 1

Rule 2

Rule 1

Rule 1

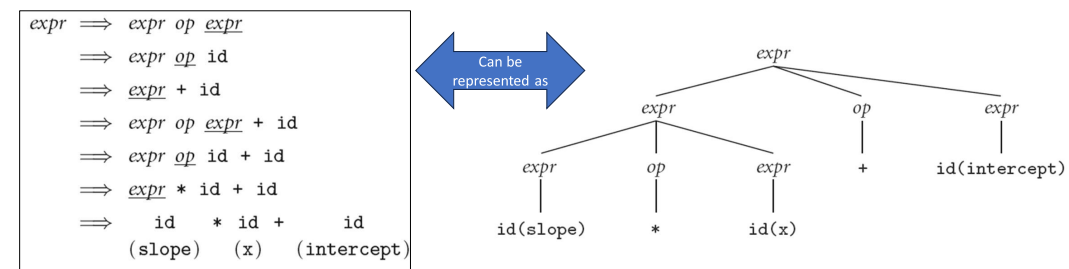
Rule 2

Rule 1

11

Parse tree

- Parse tree is a graphical representation of the derivations



12

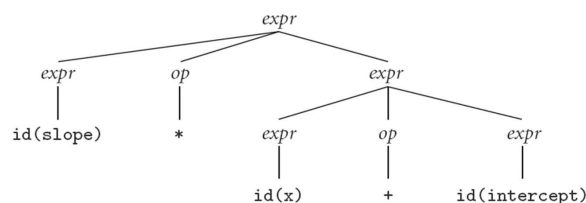
Alternative derivation

$$\begin{aligned} \text{expr} &\rightarrow \text{id} \mid \text{number} \mid - \text{expr} \mid (\text{expr}) \\ &\quad \mid \text{expr op expr} \\ \text{op} &\rightarrow + \mid - \mid * \mid / \end{aligned}$$

- Alternative parse tree for "slope * x + intercept" if we rewrite the left-most term
- Is this derivation correct ?

```

expr => expr op expr
=> id op expr
=> id * expr
=> id * expr op expr
=> id * id op expr
=> id * id + expr
=> id * id + id
(slope) (x) (intercept)
    
```



Note: Grammars which produce more than one parse tree is ambiguous – more on this in later slides

13

Another CFG with precedence

- A better version of grammar can capture **precedence**

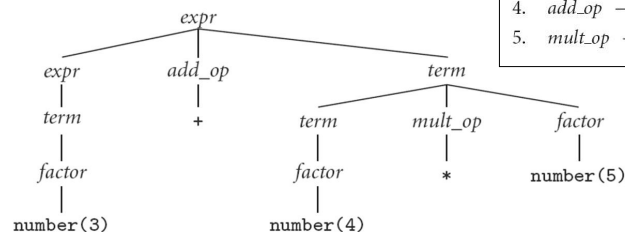
1. $\text{expr} \rightarrow \text{term} \mid \text{expr add_op term}$
2. $\text{term} \rightarrow \text{factor} \mid \text{term mult_op factor}$
3. $\text{factor} \rightarrow \text{id} \mid \text{number} \mid - \text{factor} \mid (\text{expr})$
4. $\text{add_op} \rightarrow + \mid -$
5. $\text{mult_op} \rightarrow * \mid /$

14

Derivation using the new grammars

- Parse tree for expression grammar (with left associativity) for **3 + 4 * 5**

1. $\text{expr} \rightarrow \text{term} \mid \text{expr add_op term}$
2. $\text{term} \rightarrow \text{factor} \mid \text{term mult_op factor}$
3. $\text{factor} \rightarrow \text{id} \mid \text{number} \mid - \text{factor} \mid (\text{expr})$
4. $\text{add_op} \rightarrow + \mid -$
5. $\text{mult_op} \rightarrow * \mid /$



15

Programming Language Parser

- A parser is to analyze the relationship of each word in a sentence according to the principles of grammar language
 - CFG is a generator for a context-free language (CFL)
- The parser accomplishes two tasks:
 - parsing the code
 - looking for (**syntax**) errors and generating a parse tree as the output of the phase
- Parsers are expected to parse the whole code even if some (**semantic**) errors exist in the program

16

Parsing process

- Derivation
 - a sequence of production rules, in order to get the input string.
- During parsing, we take two decisions for some sentential form of input:
 - deciding the non-terminal which is to be replaced
 - deciding the production rule, by which, the non-terminal will be replaced.
- To decide which non-terminal to be replaced with production rule, we can have two options:
 - **Left-most Derivation**
 - **Right-most Derivation**

17

Two approaches to Parsing

- Left-to-right Left-most Derivation (LL Parsing): also called 'top-down', or 'predictive' parsers
 - **scan and replace** the input with **production rules, from left to right.**
 - The sentential form derived by the left-most derivation is called the **left-sentential form.**
- Left-to-right Right-most Derivation (LR Parsing): also called 'bottom-up', or 'shift-reduce' parsers
 - **scan and replace** the input with **production rules, from right to left.**
 - The sentential form derived from the right-most derivation is called the **right-sentential form.**

18

LL vs LR

• Example

• Production rules:

- $E \rightarrow E + E$
- $E \rightarrow E * E$
- $E \rightarrow id$

• Input string: $id + id * id$

Parse Tree



• The left-most derivation is:

- $E \rightarrow E + E$
- $E \rightarrow id + E$
- $E \rightarrow id + E * E$
- $E \rightarrow id + id * E$
- $E \rightarrow id + id * id$

• The right-most derivation is:

- $E \rightarrow E + E$
- $E \rightarrow E + E * E$
- $E \rightarrow E + E * id$
- $E \rightarrow E + id * id$
- $E \rightarrow id + id * id$

19

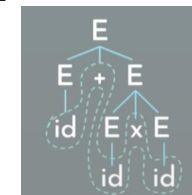
Ambiguity in CFGs

$$E \rightarrow E + E \mid E * E \mid id$$

• The left-most & right-most derivation [1]:

- $E \rightarrow E + E$
- $E \rightarrow id + E$
- $E \rightarrow id + E * E$
- $E \rightarrow id + id * E$
- $E \rightarrow id + id * id$

- $E \rightarrow E + E$
- $E \rightarrow E + E * E$
- $E \rightarrow E + E * id$
- $E \rightarrow E + id * id$
- $E \rightarrow id + id * id$



• The left-most & right-most derivation [2]:

- $E \rightarrow E * E$
- $E \rightarrow E + E * E$
- $E \rightarrow id + E * E$
- $E \rightarrow id + id * E$
- $E \rightarrow id + id * id$

- $E \rightarrow E * E$
- $E \rightarrow E * id$
- $E \rightarrow E + E * id$
- $E \rightarrow E + id * id$
- $E \rightarrow id + id * id$



20

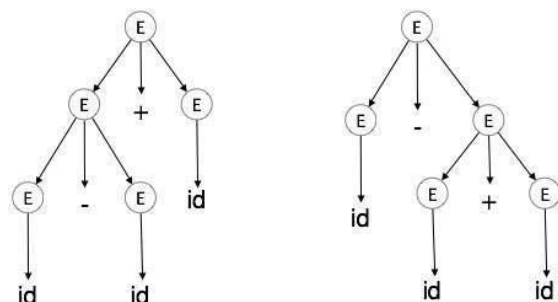
Ambiguity in Grammar

- A grammar G is said to be ambiguous if it has more than one parse tree (left or right derivation) for at least one string.

- **Example**

- $E \rightarrow E + E$
- $E \rightarrow E - E$
- $E \rightarrow \text{id}$

For the string $\text{id} + \text{id} - \text{id}$, the above grammar generates two parse trees:



https://www.tutorialspoint.com/compiler_design/compiler_design_syntax_analysis.htm

21

Is ambiguity bad?

- The language generated by an ambiguous grammar is said to be inherently ambiguous
- Ambiguity in grammar is not good for a compiler construction
- No method can detect and remove ambiguity automatically,
 - but it can be removed by either re-writing the whole grammar without ambiguity, or
 - by setting and following associativity and precedence constraints

22

Ambiguity in associativity

- **Associativity**
 - If an operand has operators on both sides, the side on which the operator takes this operand is decided by the associativity of those operators.
 - If the operation is left-associative, then the operand will be taken by the left operator or if the operation is right-associative, the right operator will take the operand.
- **Left associative example**
 - Operations such as Addition, Multiplication, Subtraction, and Division are left associative.
 - If the expression contains: id op id op id , it will be evaluated as: $(\text{id op id}) \text{ op id}$
 - For example, $(\text{id} + \text{id}) + \text{id}$
- **Right associative example**
 - Operations like Exponentiation are right associative, i.e., the order of evaluation in the same expression will be: $\text{id op} (\text{id op id})$
 - For example, $\text{id} \wedge (\text{id} \wedge \text{id})$

23

Ambiguity in precedence

- If two different operators share a common operand, the precedence of operators decides which will take the operand.
 - That is, $2+3*4$ can have two different parse trees,
 - one corresponding to $(2+3)*4$ and
 - another corresponding to $2+(3*4)$.
 - By setting precedence among operators, this problem can be easily removed.
 - As in the previous example,
 - mathematically $*$ (multiplication) has precedence over $+$ (addition),
 - so the expression $2+3*4$ will always be interpreted as: $2 + (3 * 4)$
- These methods decrease the chances of ambiguity in a language or its grammar

24

Real-world example: Python's grammars

```
# If statement
# -----

if_stmt:
    | 'if' name_expression ':' block elif_stmt
    | 'if' name_expression ':' block [else_block]

elif_stmt:
    | 'elif' name_expression ':' block elif_stmt
    | 'elif' name_expression ':' block [else_block]
else_block:
    | 'else' ':' block
```

<https://docs.python.org/3/reference/grammar.html>

25

References

- Scott, M. L. (2016). Programming Language Pragmatics.
- Tutorialspoint, <https://www.tutorialspoint.com/>, accessed on 11/22/2023
- <https://docs.python.org/3/reference/grammar.html>

26