

# Syntax Analysis: Scanning and Parsing

204315: OPL

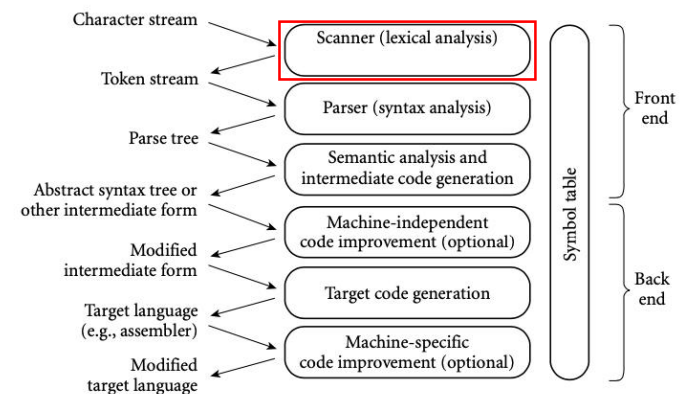
## Outline

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

2

## Scanner

## Compilation Process Overview (Recap)



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

4

## Scanner (Recap)

- Program tokenization
- Uses **regular expression**
  - Regular expressions have the capability to express languages.

```
int main() {  
    int i = getint(), j = getint();  
    while (i != j) {  
        if (i > j) i = i - j;  
        else j = j - i;  
    }  
    putint(i);  
}
```



```
int      main      (      )      {      int      i      =  
getint   (      )      ,      j      =      getint   (  
)      ;      while (      i      !=      j      )  
{      if      (      i      >      j      )      i  
=      i      -      j      ;      else      j      =  
j      -      i      ;      }      putint (      i
```

5

## Scanner

- responsible for
  - tokenizing source code
  - removing comments
  - (often) dealing with *pragmas* (i.e., compiler directives)
  - saving text of identifiers, numbers, strings
  - saving source locations (file, line, column) for error messages

6

## Scanner/Lexical Analysis

- Unlike natural languages such as English or Thai, **computer languages must be precise.**
  - To provide the precision, **language designers use formal syntactic (syntax) and semantic notation.**
- Different programming languages
  - often provide features with very **similar semantics** but very **different syntax**.
  - It is generally much easier to learn a new language if one is able to identify the common semantic ideas beneath the unfamiliar syntax.

7

## Scanner/Lexical Analysis

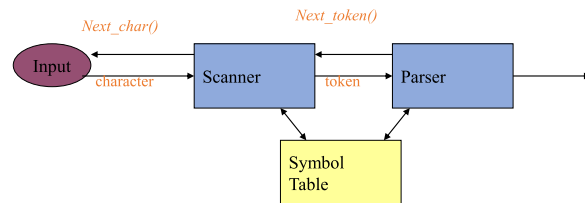
- **Lexical analysis** is the first phase of a compilation process.
  - reads/scans character streams from the source code
  - checks for legal tokens
  - passes the data to the syntax analyzer when it demands.
- **Lexical analyzer** needs to scan and identify only a finite set of valid string/token
  - It searches for the pattern defined by the language rules.
- **Regular expressions** have the capability to express finite languages by defining a pattern for finite strings of symbols.
  - The grammar defined by regular expressions is known as **regular grammar**
  - The language defined by regular grammar is known as **regular language**.

**Syntax analyzer:** checks if the source code follows the grammatical rules of the programming language.

8

## Scanner input/output

- INPUT: sequence of characters
- OUTPUT: sequence of tokens



9

## Some Definitions

- **token** – set of strings that is meaningful in source language
- **pattern** – a rule describing a set of string
- **lexeme** – a sequence of characters that matches the patterns of tokens

10

## Some Examples

Token	Pattern	Sample Lexeme
while	while	while
relation_op	=   !=   <   >	<
integer	(0-9)*	42
string	Characters between “ ”	“hello”

11

## Regular Expression (RE)

- A way to describe pattern that specifies how tokens are generated
- Examples of RE for : *number token, integer token, real token, ...*

$number \rightarrow integer \mid real$

$integer \rightarrow digit \, digit^*$

$real \rightarrow integer \, exponent \mid decimal \, (exponent \mid \epsilon)$

$decimal \rightarrow digit^* \, ( \, . \, digit \mid digit \, . ) \, digit^*$

$exponent \rightarrow (e \mid E) \, ( + \mid - \mid \epsilon ) \, integer$

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

12

## RE and their operations

- RE is one of the followings
  - A character
  - The empty string, denoted by  $\epsilon$
  - Two regular expressions **concatenated**
  - Two regular expressions separated by **|** (i.e., or)
  - A regular expression followed by the **Kleene star \*** (concatenation of zero or more strings)

$number \rightarrow integer \mid real$   
 $integer \rightarrow digit \, digit^*$   
 $real \rightarrow integer \, exponent \mid decimal \, (exponent \mid \epsilon)$   
 $decimal \rightarrow digit^* \, ( \, digit \mid digit \, . \, ) \, digit^*$   
 $exponent \rightarrow (e \mid E) \, ( + \mid - \mid \epsilon ) \, integer$   
 $digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

13

## Check your understanding

- Are these lexemes generated by these REs?

• 3.1415	$number \rightarrow integer \mid real$
• 1.1.0	$integer \rightarrow digit \, digit^*$
• 15e-5	$real \rightarrow integer \, exponent \mid decimal \, (exponent \mid \epsilon)$
• 3E++	$decimal \rightarrow digit^* \, ( \, digit \mid digit \, . \, ) \, digit^*$
• .357	$exponent \rightarrow (e \mid E) \, ( + \mid - \mid \epsilon ) \, integer$
	$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

14

## Recognizing Regular Expression

- A mechanical approach to check if a string is generated by a particular RE is by constructing and using a Deterministic Finite Automaton (DFA)
- A DFA is a [finite-state machine](#) that accepts or rejects a given [string](#) of symbols, by running through a state sequence uniquely determined by the string
- Multiple version of DFAs can be constructed from a single RE

15

## Deterministic Finite Automaton

- A deterministic finite automaton M is a 5-tuple,  $(Q, \Sigma, \delta, q_0, F)$  consisting of
  - a finite set of states Q
  - a finite set of input symbols called the alphabet  $\Sigma$
  - a transition function  $\delta : Q \times \Sigma \rightarrow Q$
  - an initial or start state  $q_0 \in Q$
  - a set of accept states  $F \subseteq Q$

16

## DFA Example

- Given RE =  $(1^*) (0 (1^*) 0 (1^*))^*$

- Possible valid strings

- 10101
- 0101
- 0000

- Invalid string

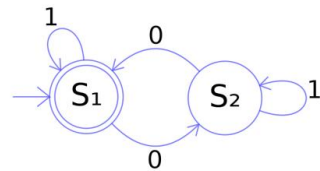
- 000
- 011

- The corresponding DFA is

$M = (Q, \Sigma, \delta, q_0, F)$  where

- $Q = \{S_1, S_2\}$
- $\Sigma = \{0, 1\}$
- $q_0 = S_1$
- $F = \{S_1\}$  and
- $\delta$  is defined by the following state transition table:

	0	1
$S_1$	$S_2$	$S_1$
$S_2$	$S_1$	$S_2$

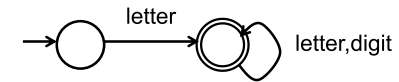


[https://en.wikipedia.org/wiki/Deterministic\\_finite\\_automaton](https://en.wikipedia.org/wiki/Deterministic_finite_automaton)

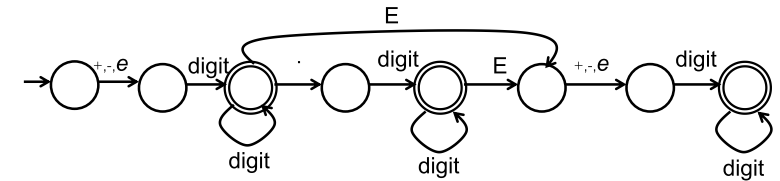
17

## More DFA examples

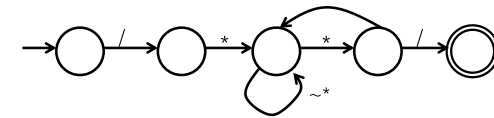
- Identifier



- Numeric



- C comment



18  
Ref: 2301373 Introduction to Compilers

## Remarks

- We run the machine (DFA) over and over to get one token after another
  - always take the longest possible token from the input  
thus foobar is foobar and never f or foo or foob
  - more to the point, 3.14159 is a real const and never 3, ., and 14159
- Regular expressions "generate" a regular language; DFAs "recognize" it

19

## How to construct a Scanner?

- Manual
  - Write down regular expressions for your language syntax
  - Construct a DFA for the regular expressions
  - Write a program (as a set of nested IFs) which walks the DFA
- Semi-automatic
  - Write down regular expressions for your language syntax
  - Use automatic softwares such as *lex*, *flex*, *jflex* that reads REs and source code and produce the tokenization result

20

## References

- PPT of Lecture 2: Lexical Analysis, CS 540, George Mason University
- PPT of 2301373 Introduction to Compilers
- Scott, M. L. (2016). Programming Language Pragmatics.

21

## What's next?

- **Parsing**
  - LL Parsing
  - LR Parsing

22