

CSCI 2320

Principles of Programming Languages

**Lexical Analysis**

Reading: Chapter 3 of Tucker-Noonan

---

MOHAMMAD T. IRFAN



# Plan

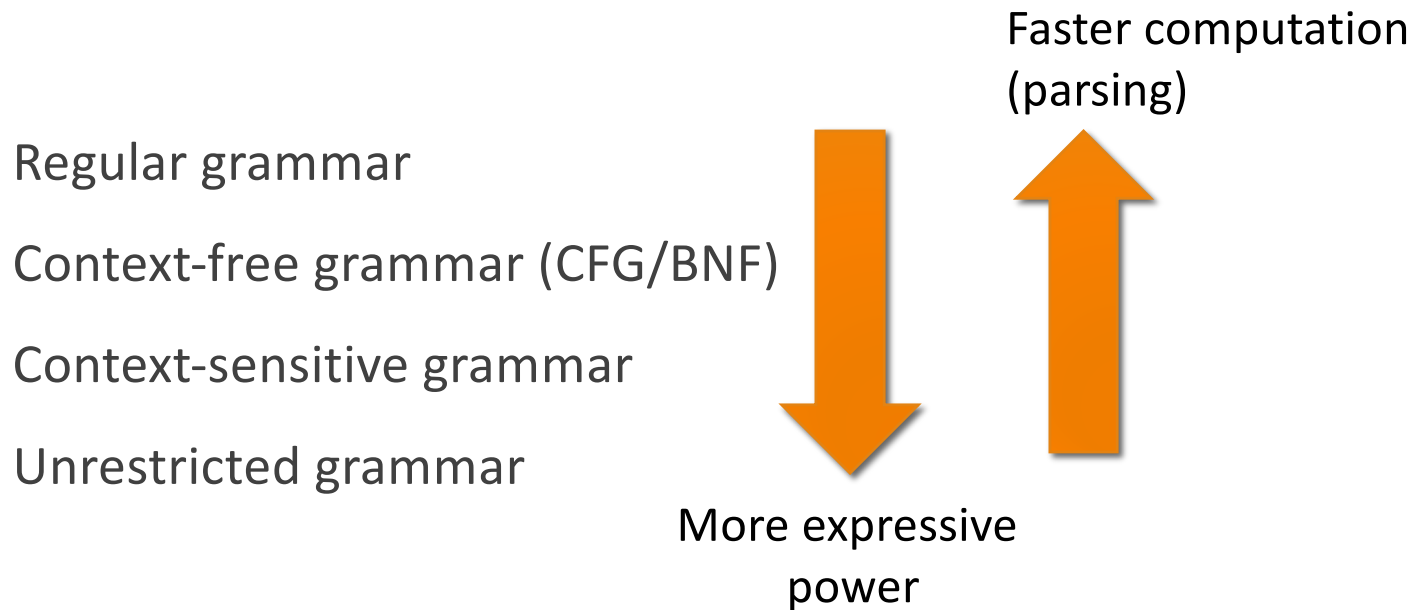
---

Chomsky's Hierarchy

Lexical Analysis

# Chomsky's Hierarchy

---



Let's dive into these

# Lexical Analysis

---



# Lexical Analysis

---

Input: Lexeme(s) = **sequence of input characters having a collective meaning**

Output: Tokens (representing lexemes)

Discard: whitespace, comments

```
int count = 10;
```

Lexemes

int

count

=

10

;

Tokens

Type

Ident  
ifier

=

IntLi  
teral

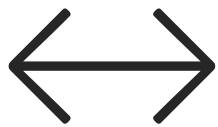
;

# Why do lexical analysis separately?

---

- 75% of compiling time spent in lexical analysis
- Simpler, faster grammar for parsing
  - Next: how?
- Take care of OS idiosyncrasies

# Regular Expressions (regex)



# Regular Grammar

---

CONSTRUCTIVE PROOF LATER ...

# Recursive definition of a regex

## Regex Meaning

$x$	A character $x$ from an alphabet is a regex	(base case)
$\epsilon$	The empty string $\epsilon$ is a regex	(base case)
$MN$	<i>Concatenation</i> : $M$ followed by $N$ is a regex	if $M$ and $N$ are regex
$M N$	<i>Alteration</i> : $M$ or $N$ is a regex	if $M$ and $N$ are regex
$M^*$	<i>Kleen star</i> : 0 or more occurrences of $M$ is a regex	if $M$ is a regex
$M^+$	<i>Kleen plus</i> : 1 or more occurrences of $M$ is a regex	if $M$ is a regex
$M?$	0 or 1 occurrence of $M$ is a regex	if $M$ is a regex

# Regex meta symbols

[set of char]      Take exactly 1 character from a set of characters.

[aeiou] means a | e | i | o | u (that is, exactly one of the five vowels).

[0-9] means 0 | 1 | ... | 9. Here, - is the range meta symbol.

{Z}      A reference to a regex labeled Z (note the difference with EBNF's { }).

( ... )      Grouping and additionally, in regex libraries, capturing.

.

Means any single character.

^

Means not (or excluding); e.g., [a-zA-Z^aeiouAEIOU] means a consonant.

\

Escaped character; e.g., \. means dot, \( means the literal left parenthesis and not the meta symbol for grouping.

C Lite

Regular Definitions

---



# CLite regular definition

---

## Category

## Definition

AnyChar

[ -~]

From space (ASCII 27) to tilde (126)

Letter

[a-zA-Z]

Digit

[0-9]

Whitespace

[ \t]

Space and tab

Eol

\n

## Category

Keyword

Identifier

IntegerLit

FloatLit

CharLit

## Definition

bool | char | else | false | float |

if | int | main | true | while

{Letter}({Letter} | {Digit})\*

{Digit}+

{Digit}+\.{Digit}+

'{AnyChar}'

## Category

## Definition

Operator

= | || | && | == | != | < |  
<= | > | + | - | \* | / | ! | [ | ]

Separator

; | . | { | } | ( | )

Comment

// ({AnyChar} | {Whitespace})\* {Eol}

CFG  
vs  
Regular grammar  
vs  
Regex

---



# Describe the language:

---

1.  $0(0|1)^+0$


2.  $((\epsilon|0)1^*)^*$

3.  $0^*10^*10^*10^*$

4.  $(00|11)^*$

# Write regular expression for:

---

1. All strings of lowercase letters, where letters appear in ascending order.
  2. All strings of letters containing vowels in order.
  3. Strings of digits (0, 1, 2) with no consecutively repeated digit.
- 

# Finite State Automata (FSA) --- DFA, NFA

---

BEHIND THE SCENE OF REGULAR EXPRESSIONS

READING:  
TEXTBOOK AND HANDOUT (SCOTT)

# Known algorithms

---

1. Regular expression  $\rightarrow$  NFA
2. NFA  $\rightarrow$  DFA

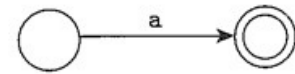
Language designer  $\rightarrow$   
implementation (parsing)

3. DFA  $\rightarrow$  regular expression

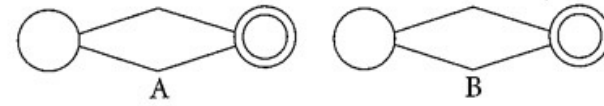
State elimination

Regex  $\rightarrow$  NFA  $\rightarrow$  DFA  $\rightarrow$  Regex  
All 3 are equivalent!

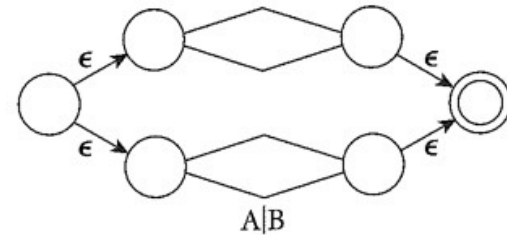
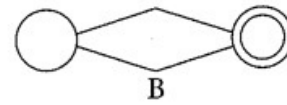
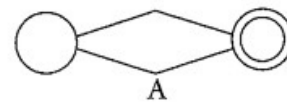
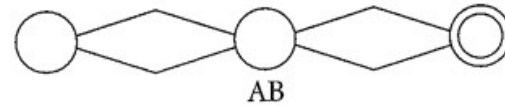
# Regex $\rightarrow$ NFA



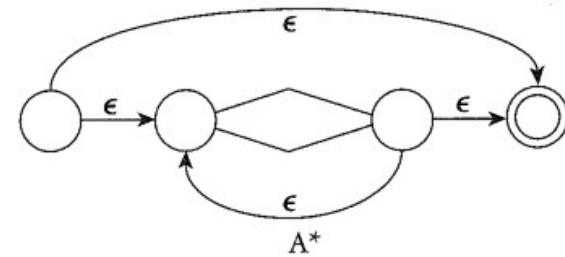
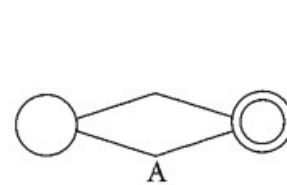
(a)



(b)



(c)



(d)

NFA  $\rightarrow$  DFA

---

Tabular method



# Regex $\leftrightarrow$ Regular Grammar

---

Regex  $\rightarrow$  NFA (2 slides ago)  
 $\rightarrow$  Regular grammar:

1. NFA state transition  $Q_i \rightarrow Q_j$  on input  $a$   
 $\Rightarrow$  Production rule  $Q_i \rightarrow a Q_j$
2. NFA final states  $F$   
 $\Rightarrow$  Production rule  $F \rightarrow \epsilon$

# Regex $\leftrightarrow$ Regular Grammar

---

Regular grammar

- NFA (state transitions follow production rules)
- DFA (tabular method)
- Regex (state elimination algorithm)

# Lexical Analysis Using Python

---

Project 1



# Python's re package

---

<https://docs.python.org/3/library/re.html>

```
import re #regex  
re.split(...) #Use regex argument to split a string into parts
```

Common string-matching regex:

## Symbol

## Definition

\d

[0-9]

\D

[^0-9]

\w

[a-zA-Z0-9\_]

\W

[^a-zA-Z0-9\_]

# Resources

---

Tucker-Noonan book chapter

Scott's book chapter (handout on Canvas)

Video lecture on NFA  $\rightarrow$  DFA: <https://youtu.be/qfVTXwuxEOY>