

Computing concepts for Compiler Design

Design of Lexical Analyzer

Concepts and Notations

- **Set:** An unordered collection of unique elements

$S_1 = \{ a, b, c \}$ $S_2 = \{ 0, 1, \dots, 19 \}$ empty set: \emptyset

membership: $x \in S$ union: $S_1 \cup S_2 = \{ a, b, c, 0, 1, \dots, 19 \}$

universe of discourse: U subset: $S_1 \subset U$

complement: if $U = \{ a, b, \dots, z \}$, then $S_1' = \{ d, e, \dots, z \} = U - S_1$

- **Alphabet:** A finite set of symbols

- Examples:

- Character sets: [ASCII](#), [ISO-8859-1](#), Unicode

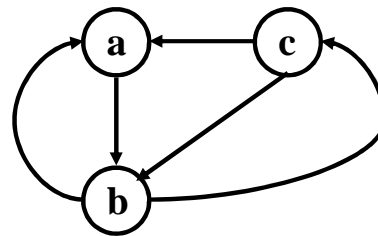
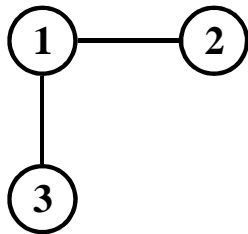
- $S_1 = \{ a, b \}$ $S_2 = \{ \text{Spring, Summer, Autumn, Winter} \}$

- **String:** A sequence of zero or more symbols from an alphabet

- The empty string: ϵ

Concepts and Notations

- **Language**: A set of strings over an alphabet
 - Also known as a **formal language**; may not bear any resemblance to a **natural language**, but could model a subset of one.
 - The language comprising **all** strings over an alphabet Σ is written as: Σ^*
- **Graph**: A set of **nodes** (or **vertices**), some or all of which may be connected by **edges**.
 - An example: – A **directed graph** example:



Regular Expression

- It is a tool to express language in the form of expression.
- RE uses primitive operators for expressing language.
- The three operators used are: Union, Concatenation, Kleene Star.
- Rules for Regular Expression:
 - Every letter of alphabet set represents regular expression.
 - ϕ is regular expression.
 - ε is regular expression.
 - If r_1 and r_2 are regular expressions then (r_1) , (r_2) , $(r_1 . r_2)$ (r_1+r_2) (r_1^*) (r_2^*) will be denoted as regular expression.
 - Nothing else is regular expression.

Examples on regular expressions

| Language | Regular Expression | Remarks |
|---|--|---|
| $L = \{a^n b^m \mid n+m \text{ is even}\}$ | $(aa)^*(bb)^* + (aa)^*a(bb)^*b$ | Both n and m are even Both n and m are odd |
| $L = \{\text{All set of strings of } (0,1) \text{ with even 0s followed by odd 1s}\}$ $L = \{\text{All set of strings of } (0,1) \text{ with even 0s or odd 1s}\}$ | $(00)^* (11)^* 1$ $(00)^* + (11)^* 1$ | |
| $L = \{a,b\}$ No two consecutive letters are same | Two possibilities: $b (ab)^* a$ $a (ba)^* b$ | |
| $L = \{W \in (a,b) \text{ with length}(w)=\text{even}\}$ | $((aa) (ab) (ba) (bb))^*$ | |
| $L = \{W \in (a,b)^*, \text{ such that } W \text{ contains even number of } a\text{'s and odd number of } b\text{'s and}\}$ | $L = L1 \text{ UNION } L2$ AND if L1 and L2 are regular, then L is regular. | |

Examples: Regular Expressions

1. 0^*10^* , $L(0^*10^*) = \{w \mid w \text{ contains exactly a single } 1\}$
2. $\Sigma^*1\Sigma^*$, $L(\Sigma^*1\Sigma^*) = \{w \mid w \text{ contains at least one } 1\}$
3. $\Sigma^*001\Sigma^*$, $L(\Sigma^*001\Sigma^*) = \{w \mid w \text{ contains } 001 \text{ as a substring}\}$
4. $(\Sigma\Sigma)^*$, $L(\Sigma\Sigma)^* = \{w \mid w \text{ is a string of even length}\}$
5. $(\Sigma\Sigma\Sigma)^*$, $L(\Sigma\Sigma\Sigma)^* = \{w \mid \text{the length of } w \text{ is a multiple of three}\}$
6. $01 \cup 01$, $L(01 \cup 01) = \{01, 01\}$
7. $0\Sigma^*0 \cup 1\Sigma^*1 \cup 0 \cup 1$, $L(0\Sigma^*0 \cup 1\Sigma^*1 \cup 0 \cup 1) = \{w \mid w \text{ starts and ends with the same symbol}\}$
8. $(0 \cup \epsilon)1^*$, $L((0 \cup \epsilon)1^*) = \{01^* \cup 1^*\}$
9. $(0 \cup \epsilon)(1 \cup \epsilon)$, $L((0 \cup \epsilon)(1 \cup \epsilon)) = \{\epsilon, 0, 1, 01\}$

Examples:

Examples of Regular Languages

- $L = \{x \in \{0,1\}^* \mid |x| \text{ is even}\}$
 - Any string of even length can be obtained by concatenating strings length 2.
 - Any concatenation of strings of length 2 will be even
 - $L = \{00, 01, 10, 11\}^*$
- Regular expressions describing L:
 - $(00 + 01 + 10 + 11)^*$
 - $((0 + 1)(0 + 1))^*$

Examples of Regular Languages

- $L = \{x \in \{0,1\}^* \mid x \text{ does not end in } 01 \}$
 - If x does not end in 01, then either
 - $|x| < 2$ or
 - x ends in 00, 10, or 11
 - A regular expression that describes L is:
 - $\Lambda + 0 + 1 + (0 + 1)^*(00 + 10 + 11)$

Examples of Regular Languages

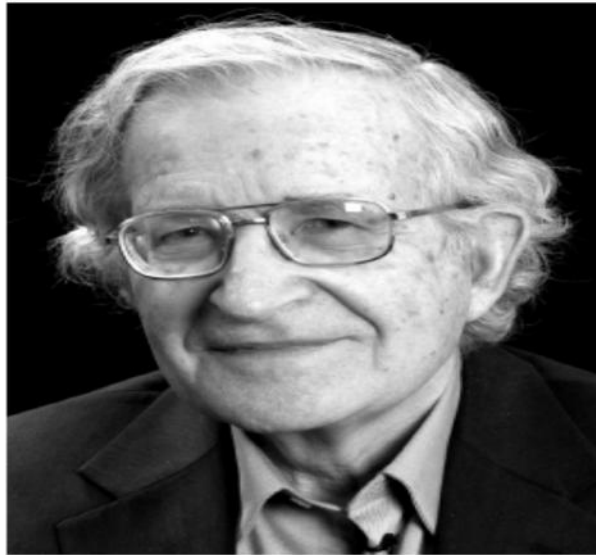
- $L = \{x \in \{0,1\}^* \mid x \text{ contains an odd number of 0s}\}$
 - Express $x = yz$
 - y is a string of the form $y = 1^i 0 1^j$
 - In z, there must be an even number of additional 0s or $z = (01^k 0 1^m)^*$
 - x can be described by $(1^* 0 1^*)(01^* 0 1^*)^*$

Examples

- ▶ Let $\Sigma = \{a, b\}$
 - ▶ $a|b$ denotes $\{a, b\}$
 - ▶ $(a|b)(a|b)$ denotes $\{aa, ab, ba, bb\}$
i.e., $(a|b)(a|b) = aa|ab|ba|bb$
 - ▶ a^* denotes $\{\epsilon, a, aa, aaa, \dots\}$
 - ▶ $(a|b)^*$ denotes the set of all strings of a 's and b 's (including ϵ)
i.e., $(a|b)^* = (a^*b^*)^*$
 - ▶ $a|a^*b$ denotes $\{a, b, ab, aab, aaab, aaaab, \dots\}$

Context Free Grammars

Context-Free Grammars



Noam Chomsky
(linguist, philosopher, logician, and activist)

*“ A **grammar** can be regarded as a **device** that enumerates the **sentences** of a **language**. We study a sequence of restrictions that limit grammars first to **Turing machines**, then to two types of systems from which a phrase structure description of a generated language can be drawn, and finally to finite state Markov sources (**finite automata**). ”*

Chomsky Hierarchy of Grammar

| Description | Type | Terminology |
|---------------------------------------|---|--|
| Un-restricted Grammar | Type – 0 | No restriction on writing production rules. |
| Context Sensitive Grammar | Type – 1 $aXb \rightarrow aYb$ If X = surrounded by "a" and "b" replace X with Y | Depending on context, non-terminals are expanded. Tag rules of HTML language |
| Context Free Grammar | Type -2 | The production rule will only have one non-terminal on LHS Permitted rules $A \rightarrow a$ $A \rightarrow B$ $A \rightarrow$ any combination of terminal and non-terminals $A \rightarrow \epsilon$ |
| Regular Grammar Restricted grammar | Type - 3 | Permitted rules: $A \rightarrow aB$ or $A \rightarrow a$ or $A \rightarrow \epsilon$ |

Formal Definition of Context Free Grammar

- CFG is defined as: $G = (V, T, P, S)$
 - V = set of non-terminals
 - T = set of terminals
 - P = Production rules: of form $A \rightarrow \gamma$ where $\gamma = (VUT)^*$
 - γ can be ϵ or single terminal or single non-terminal symbol.
- γ can be any combination of terminal and non-terminals
- Why CFG:
- Since there is no restriction on right side of production rule, it is denoted as context free, i.e., ir-respective of surrounding, non-terminal symbol can be replaced by RHS if there exist production rule.

Design of CFG

- Example: 1: Balanced parenthesis grammar
 - $S \rightarrow (S)$
 - $S \rightarrow SS$
 - $S \rightarrow \varepsilon$

- Example 2: Grammar for $L = \{a^n b^n \mid n \geq 0\}$
 - $S \rightarrow aSb$
 - $S \rightarrow \varepsilon$

Design of CFG

- Example: 3: Even length palindrome
 - $S \rightarrow aSa$
 - $S \rightarrow bSb$
 - $S \rightarrow \varepsilon$

- Example 4: Grammar for $L = \text{Equal number of } a\text{'s and } b\text{'}$
 - $S \rightarrow aSb$
 - $S \rightarrow bSa$
 - $S \rightarrow SS$
 - $S \rightarrow \varepsilon$

Design of CFG

- Example: 4: $L = \{W \in (0,1)^* \mid W \text{ contains three } 1\text{'s}\}$
 - $S \rightarrow X1X1X1$
 - $X \rightarrow 0X \mid 1X \mid \epsilon$
- Example 5: $L = \{W \in (0,1)^* \mid W \text{ length is odd and middle symbol is } 0\}$
 - $S \rightarrow 0S0 \mid 1S1 \mid 0S1 \mid 1S0$
- Example 6: $L = \{a^i b^j c^k \mid i+j = k\}$
 - $S \rightarrow aSc \mid X$
 - $X \rightarrow bXc \mid \epsilon$

Write CFGs for the following languages:

1. Strings ending with a 0
2. Strings containing even number of 1's
3. palindromes over $\{0, 1\}$
4. $L = \{a^i b^j : i \leq 2j\}$ or $L = \{a^i b^j : i < 2j\}$ or $L = \{a^i b^j : i \neq 2j\}$
5. $L = \{a^i b^j c^k : i = k\}$
6. $L = \{a^i b^j c^k : i = j\}$
7. $L = \{a^i b^j c^k : i = j + k\}$.
8. $L = \{w \in \{0, 1\}^* : |w|_a = |w|_b\}$.

Reading assignment

- Read:
 - Cross Compiler
 - Example of cross compiler
 - Bootstrapping
 - Example of bootstrapping