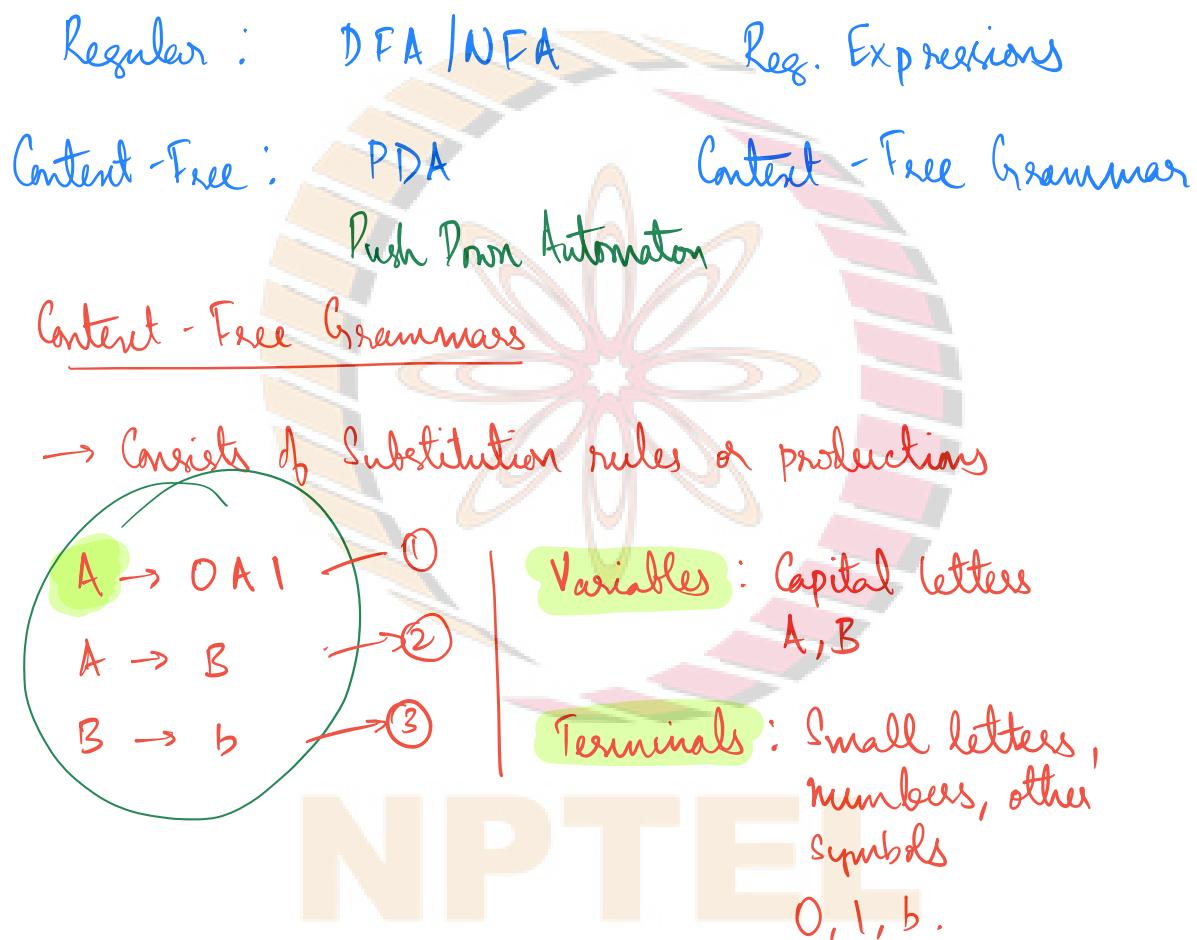


Content-Free languages

- Next step from regular languages
- Useful in parsing programming languages.



Example derivations : $A \rightarrow B \rightarrow \underline{b}$

$$A \rightarrow OA1 \rightarrow OBI \rightarrow \underline{Ob1}$$

$$A \rightarrow \underline{OA1} \rightarrow O \underline{OA1} I \rightarrow OOBI \rightarrow \underline{OObII}$$

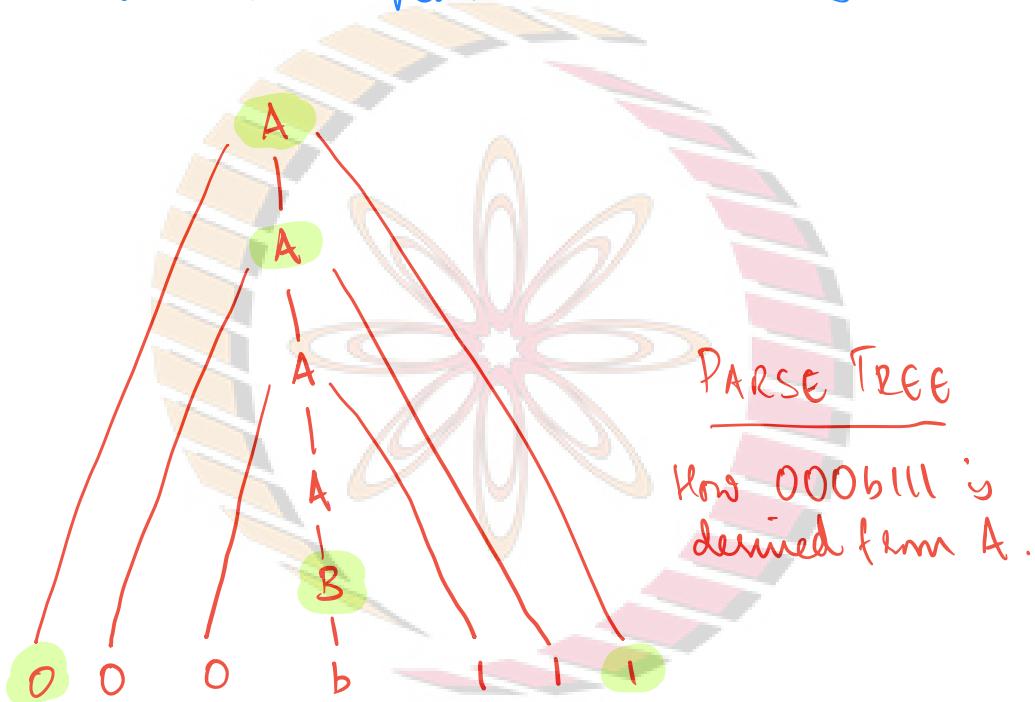
$$A \rightarrow \underline{OA1} \rightarrow OOAII \rightarrow OOOAIII \rightarrow OOOBIII \rightarrow \underline{OOObIII}$$

..

Strings generated : b, 0b1, 00b11, 000b111, ...

Start variable : On the LHS of the first rule.

→ Start from start variable, repeatedly replace the variables in the resulting strings using the rules. Repeat till no variables remain.



If we write a program, the compiler needs to check for properly nested parentheses.

$$L = \{ w \in \{ [,]\}^* \mid w \text{ is properly nested} \}$$

Examples: $\underline{[}[\underline{[} \underline{]}]] \underline{1}, \underline{[} \underline{[} \underline{]} \underline{]} \underline{[} \underline{]}, \underline{[} \underline{[} \underline{[} \underline{]} \underline{]} \underline{]} \underline{[} \underline{]} \underline{]} \underline{1}$

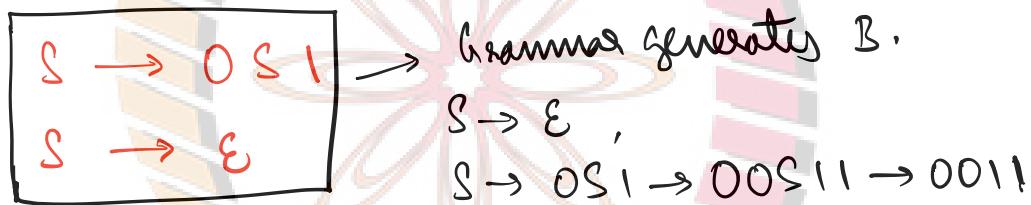
$$\Sigma = \{ [,]\}$$

$$A = [^*]^* \rightarrow A \text{ is regular.}$$

$$L \cap A = \{ [^n]^n \mid n \geq 0 \} \quad \text{is not regular}$$

Hence L is not regular. We will see that L is content-free.

$$B = \{ 0^n 1^n \mid n \geq 0 \} \rightarrow \text{Similar to } L \cap A$$



Even natural languages like English follow such a grammar. We have words and we generate sentences in various ways from NOUN, VERB, PREPOSITION, ADJECTIVE etc. (read example in the book, page 101).

Examples: (1) $S \rightarrow (A)$

$$A \rightarrow \underbrace{\epsilon \mid aA \mid ASA}_{A \text{ can generate } \epsilon, \text{ or } aA, \text{ or } ASA}$$

Variables: S, A

Terminals: $a, (,)$

$$S \rightarrow (A) \rightarrow (aA) \rightarrow (aaA) \rightarrow (aaaA)$$

- ... $S \rightarrow (A) \rightarrow (aA) \rightarrow (aaA)$

Variables: S
 Terminals: a, b
 Derives $\{a^n b^n \mid n \geq 0\}$

$(2) \quad S \rightarrow \epsilon \mid aSb$

$\rightarrow (aa\cdots a) - \text{number} \rightarrow (aa\cdots a)A \rightarrow (aa\cdots a)$

$S \rightarrow aSb \rightarrow aab \rightarrow aaaabb \rightarrow aaaaabb$
 $S \rightarrow aSb \rightarrow ab \rightarrow aabb$

Def 2.2: A **content-free grammar** (CFG) is a

4-tuple (V, Σ, R, S) , where

(1) V is a finite set, called **variables**

(2) Σ is a finite set, disjoint from V ,
called **terminals**.

(3) R is a set of rules, each of the form

Variable \rightarrow String of variables and

$A \rightarrow a, \quad A \rightarrow aSbAa$ terminals

(4) $S \in V$ is the **start variable**.

A **content-free language** (CFL) is a language that can be generated from a content-free grammar (CFG).

Yields: One step by application of a rule

\Rightarrow

$uA\omega$ \Rightarrow $u\vee\omega$ (using $A \rightarrow v$)

Defines: u derives v , $\boxed{u \xrightarrow{*} v}$ if $u=v$ or
if there is a sequence u_1, u_2, \dots, u_k such that

$$u = u_1 \Rightarrow u_2 \Rightarrow u_3 \Rightarrow \dots \Rightarrow u_k = v$$

The language of the grammar G , denoted $L(G)$,
is given by

$$L(G) = \{ w \in \Sigma^* \mid S \xrightarrow{*} w \}$$

Why context-free?

The rules and derivations are not based on the
content of the variable. Rules are not of the type

$$a \vee b \rightarrow a \times b$$

$$b \vee a \rightarrow b \vee a$$

$$\vee \rightarrow x$$

Rules do not depend on the content. There are
content-sensitive languages, though we won't see
them in this course.

Example 2.3: $h = \{\{S\}, \{a, b\}, R, S\}$

C. $[]$ $R: S \rightarrow aSb \mid SS \mid \epsilon$

If $a = [$ and $b =]$, this yields the set of all properly nested brackets.

$S \rightarrow aSb \rightarrow aaSbb \rightarrow aabb$
 $S \rightarrow SS \rightarrow aSbS \rightarrow aaSbbS \rightarrow aabbS$
 $aabbab \leftarrow aabbabS \leftarrow [] [] []$

Designing CFG's



1. If A and B are CFG's, we can construct a CFG that generates the union language by having a rule

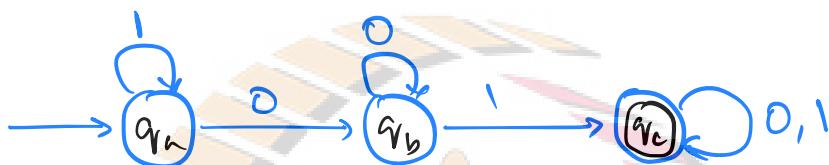
New start var. $S \rightarrow S_A \mid S_B$

CFL's are closed under union

along with the rules of A and B.

2. All regular languages can be expressed using CFG's.

- One variable for each state of the CFG.
- v_0 corresponds to q_0 (start state).
- $\delta(q_i, \alpha) = q_j$ corresponds to $v_i \rightarrow \alpha v_j$
- For $q_i \in F$, we add $v_i \rightarrow \epsilon$.



These var: A, B, C.

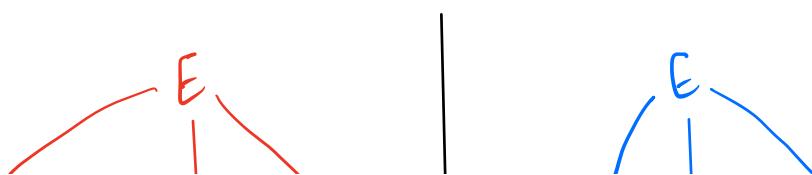
$$\begin{array}{l}
 A \rightarrow 0B \mid 1A \\
 B \rightarrow 0B \mid 1C \\
 C \rightarrow 0C \mid 1C \mid \epsilon
 \end{array}$$

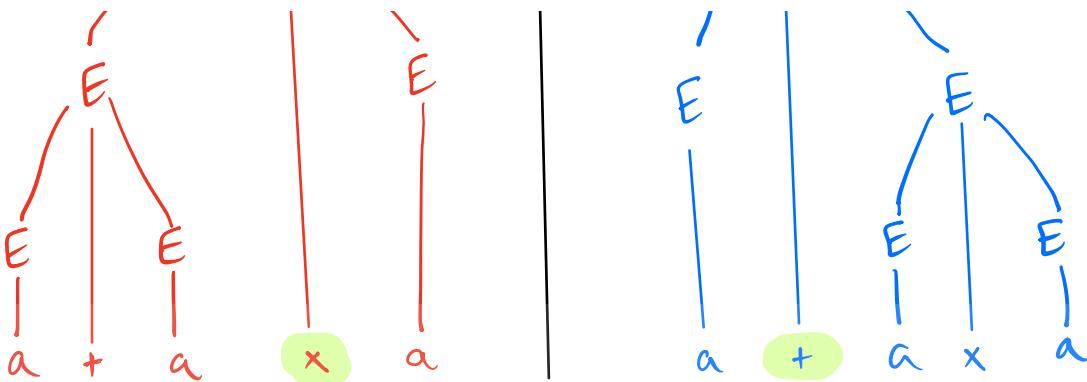
Exercise: Verify that the grammar generates the same language as the DFA.

Ambiguity

Var: E
Term: a, +, x, (,)

$$G: E \rightarrow E+E \mid E \times E \mid (E) \mid a$$





When the parser parses ' $a + a x a$ ' it is not sure how this expression was derived. Does the '+' or 'x' have precedence?

In English : She called the man with an Iphone.

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T \times F \mid F \\ F \rightarrow (E) \mid a \end{array}$$

This grammar generates the same language as above, but is not ambiguous.

Leftmost Derivation : The leftmost variable is expanded first.

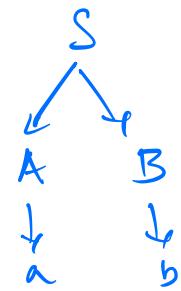
Example : $S \rightarrow A B,$

$$A \rightarrow a$$

$$B \rightarrow b$$

$$S \rightarrow AB \rightarrow aB \rightarrow ab$$

$$S \rightarrow AB \rightarrow A b \rightarrow ab$$



Def 2.7: A string w is derived ambiguously in a CFG G if it has two or more distinct leftmost derivations. Grammar G is ambiguous if it generates some string ambiguously.

NPTEL