

# CS310: Automata Theory

## Context Free Grammars

Paritosh Pandya

Indian Institute of Technology, Bombay  
Course URL: <https://cse.iitb.ac.in/~pandya58/CS310/automata.html>

Autumn, 2023

- Context free grammars and Generated languages  
Kozen L19, L20.
- Grammar Normal Forms  
Kozen L21
- Pumping Lemma for CFL  
Kozen L22
- Properties of CFL

# Phrase Structured Grammars

Noam Chomsky (Founded Computational Linguistics (1950))

- $S \rightarrow NP \ VP$   
 $NP \rightarrow (AP) \ Noun$   
 $AP \rightarrow Adj \ (AP)$   
 $VP \rightarrow Verb \ (AdvP)$   
 $AdvP \rightarrow Adv \ (AdvP)$

# Phrase Structured Grammars

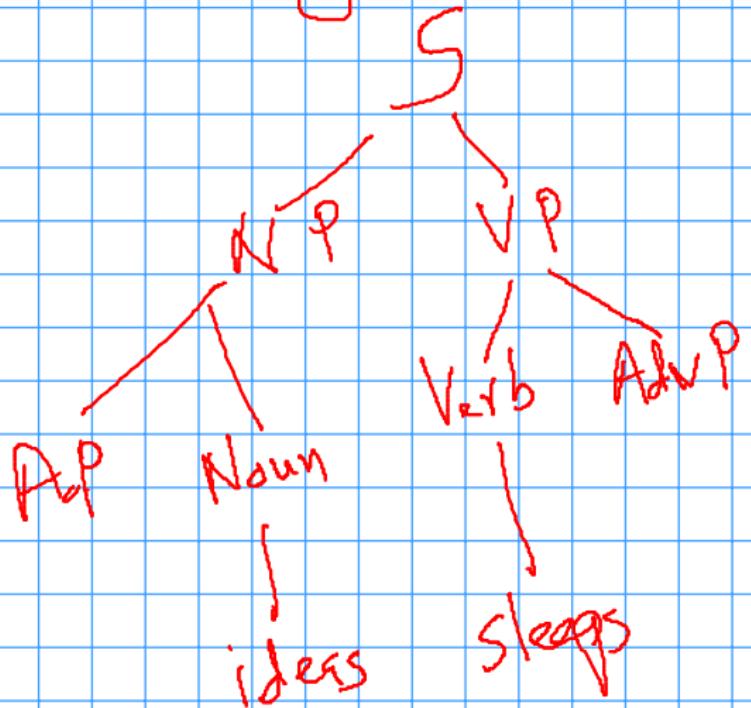
Noam Chomsky (Founded Computational Linguistics (1950))

- $S \rightarrow NP \ VP$
- $NP \rightarrow (AP) \ Noun$
- $AP \rightarrow Adj \ (AP)$
- $\underline{VP} \rightarrow Verb \ (AdvP)$
- $AdvP \rightarrow Adv \ (AdvP)$

$S \rightarrow NP \cdot VP$   
 $\rightarrow NP \cdot \underline{Verb} \cdot AdvP$   
 $\rightarrow NP \ sleep \ \underline{AdvP}$   
 $\rightarrow NP \ sleep \ \underline{Adv}$   
 $\rightarrow NP \ sleep \ furiously$   
 $\rightarrow AP \ Noun \ sleep \ furious$

- Derive English Sentences  
colorless green ideas sleep furiously

# Parsing



# Phrase Structured Grammars

Noam Chomsky (Founded Computational Linguistics (1950))

- $S \rightarrow NP \ VP$   
 $NP \rightarrow (AP) \ Noun$   
 $AP \rightarrow Adj \ (AP)$   
 $VP \rightarrow Verb \ (AdvP)$   
 $AdvP \rightarrow \underline{Adv} \ (\underline{AdvP})$

- Derive English Sentences  
colorless green ideas sleep furiously
- Grammatically correct sentences.

# Arithmetic Expressions

E.g.  $(3 + (2*4))$  Fully Parenthesized expressions made from digits.

1)  $E \rightarrow Digit$

2)  $E \rightarrow (E + E)$

3)  $E \rightarrow (E * E)$

$E \rightarrow Digit \mid (E + E) \mid (E * E)$

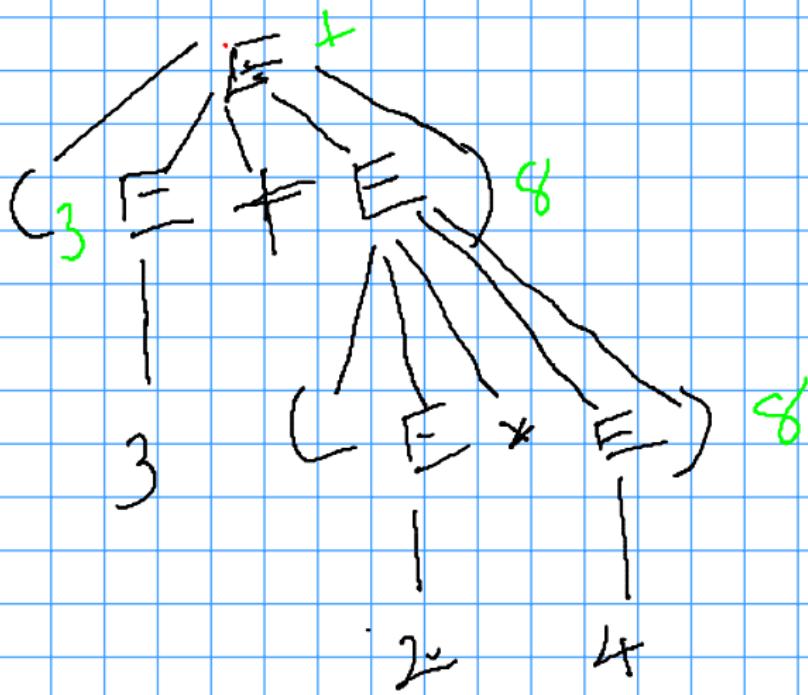
$$E \xrightarrow{1} (E + E) \xrightarrow{2} (3 + E)$$

$$\xrightarrow{3} (3 + (E * E)) \xrightarrow{\text{---}} (3 + (E * E)) \xrightarrow{\text{---}}$$

$$(3 + (2 * E)) \xrightarrow{\text{---}}$$

$$(3 + (2 * 4))$$

- Context free Grammars (CFG).
- Syntactic categories (Nonterminals) like  $E$ .
- Symbols of alphabet (Terminals) like  $+ * ) ( .$
- Derivation rules (Productions)
- Start Symbol like  $E$  or  $S$
- Useful for defining artificial languages and their processing.



# Context Free Grammar

CFG  $G = (N, \Sigma, P, S)$  where

- $N$  finite set of nonterminal symbols.
- $\Sigma$  finite set of terminal symbols (alphabet). Disjoint from  $N$
- $S$  start symbol where  $S \in N$ .

# Context Free Grammar

CFG  $G = (N, \Sigma, P, S)$  where

- $N$  finite set of nonterminal symbols.
- $\Sigma$  finite set of terminal symbols (alphabet). Disjoint from  $N$ .
- $S$  start symbol where  $S \in N$ .
- $P \subseteq_{fin} N \times (N \cup \Sigma)^*$  finite set of productions.

# Context Free Grammar

CFG  $G = (N, \Sigma, P, S)$  where

- $N$  finite set of nonterminal symbols.
- $\Sigma$  finite set of terminal symbols (alphabet). Disjoint from  $N$
- $S$  start symbol where  $S \in N$ .
- $P \subseteq_{fin} N \times (N \cup \Sigma)^*$  finite set of productions.
- Each production has the form  $(A, \alpha)$  with  $A \in N$  and  $\alpha \in (N \cup \Sigma)^*$ .

Abbreviated  $A \rightarrow \alpha$ . Also  $A \rightarrow \alpha_1 | \alpha_2 | \alpha_3$

# Context Free Grammar

CFG  $G = (N, \Sigma, P, S)$  where

- $N$  finite set of nonterminal symbols.
- $\Sigma$  finite set of terminal symbols (alphabet). Disjoint from  $N$
- $S$  start symbol where  $S \in N$ .
- $P \subseteq_{fin} N \times (N \cup \Sigma)^*$  finite set of productions.
- Each production has the form  $(A, \alpha)$  with  $A \in N$  and  $\alpha \in (N \cup \Sigma)^*$ .

Abbreviated  $A \rightarrow \alpha$ . Also  $A \rightarrow \alpha_1 | \alpha_2 | \alpha_3$

E.g.  $G = (\{S\}, \{a, b\}, P, S)$  with productions

$S \rightarrow aSb | \epsilon$ .

# Context Free Grammar

CFG  $G = (N, \Sigma, P, S)$  where

- $N$  finite set of nonterminal symbols.
- $\Sigma$  finite set of terminal symbols (alphabet). Disjoint from  $N$ .
- $S$  start symbol where  $S \in N$ .
- $P \subseteq_{fin} N \times (N \cup \Sigma)^*$  finite set of productions.
- Each production has the form  $(A, \alpha)$  with  $A \in N$  and  $\alpha \in (N \cup \Sigma)^*$ .

Abbreviated  $A \rightarrow \alpha$ . Also  $A \rightarrow \alpha_1 | \alpha_2 | \alpha_3$

E.g.  $G = (\{S\}, \{a, b\}, P, S)$  with productions

$$S \rightarrow aSb \mid \epsilon.$$

Derivation as a sequence of rewrites

$$S \rightarrow aSb \rightarrow \underline{a} \underline{a} \underline{S} \underline{b} \rightarrow a a a S b b b \rightarrow a a a b b b$$

$$L(G) = \{a^n b^n \mid n \geq 0\}$$

- Sentential forms  $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$ , Nonterminals  $A, B, C$ ,  
Terminals  $a, b, c$ .

- Sentential forms  $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$ , Nonterminals  $A, B, C$ , Terminals  $a, b, c$ .
- $\alpha \xrightarrow{1}_G \beta$  provided  $\alpha = \alpha_1 A \alpha_2$  and  $A \rightarrow \gamma \in P$  and  $\beta = \alpha_1 \gamma \alpha_2$ .

- Sentential forms  $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$ , Nonterminals  $A, B, C$ , Terminals  $a, b, c$ .
- $\alpha \xrightarrow[G]{1} \beta$  provided  $\alpha = \alpha_1 A \alpha_2$  and  $A \rightarrow \gamma \in P$  and  $\beta = \alpha_1 \gamma \alpha_2$ .
- $\alpha \xrightarrow[G]{0} \alpha$  and  
 $\alpha \xrightarrow[G]{n+1} \beta$  if  $\alpha \xrightarrow[G]{n} \gamma \xrightarrow[G]{1} \beta$  for some  $\gamma$ .

- Sentential forms  $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$ , Nonterminals  $A, B, C$ , Terminals  $a, b, c$ .
- $\alpha \xrightarrow[G]{1} \beta$  provided  $\alpha = \alpha_1 A \alpha_2$  and  $A \rightarrow \gamma \in P$  and  $\beta = \alpha_1 \gamma \alpha_2$ .
- $\alpha \xrightarrow[G]{0} \alpha$  and  
 $\alpha \xrightarrow[G]{n+1} \beta$  if  $\alpha \xrightarrow[G]{n} \gamma \xrightarrow[G]{1} \beta$  for some  $\gamma$ .
- Derivable  $\alpha \xrightarrow[G]{*} \beta$  if  $\alpha \xrightarrow[G]{n} \beta$  for some  $n \geq 0$ .

- Sentential forms  $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$ , Nonterminals  $A, B, C$ , Terminals  $a, b, c$ .
- $\alpha \xrightarrow{1}_G \beta$  provided  $\alpha = \alpha_1 A \alpha_2$  and  $A \rightarrow \gamma \in P$  and  $\beta = \alpha_1 \gamma \alpha_2$ .
- $\alpha \xrightarrow{0}_G \alpha$  and  
 $\alpha \xrightarrow{n+1}_G \beta$  if  $\alpha \xrightarrow{n}_G \gamma \xrightarrow{1}_G \beta$  for some  $\gamma$ .
- Derivable  $\alpha \xrightarrow{*}_G \beta$  if  $\alpha \xrightarrow{n}_G \beta$  for some  $n \geq 0$ .
- Language generated by CFG  $G$  contains all words over  $\Sigma$  derivable from  $S$   
$$L(G) = \{x \in \Sigma^* \mid S \xrightarrow{*}_G x\}$$

- Sentential forms  $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$ , Nonterminals  $A, B, C$ , Terminals  $a, b, c$ .
- $\alpha \xrightarrow[G]{1} \beta$  provided  $\alpha = \alpha_1 A \alpha_2$  and  $A \rightarrow \gamma \in P$  and  $\beta = \alpha_1 \gamma \alpha_2$ .
- $\alpha \xrightarrow[G]{0} \alpha$  and  
 $\alpha \xrightarrow[G]{n+1} \beta$  if  $\alpha \xrightarrow[G]{n} \gamma \xrightarrow[G]{1} \beta$  for some  $\gamma$ .
- Derivable  $\alpha \xrightarrow[G]{*} \beta$  if  $\alpha \xrightarrow[G]{n} \beta$  for some  $n \geq 0$ .
- Language generated by CFG  $G$  contains all words over  $\Sigma$  derivable from  $S$   
$$L(G) = \{x \in \Sigma^* \mid S \xrightarrow[G]{*} x\}$$
- A language  $L$  is context-free if  $L = L(G)$  for some CFG  $G$ .

## Examples

Grammar  $G_1$  with productions  $S \rightarrow aSb \mid \epsilon$

Grammar  $G_1$  with productions  $S \rightarrow aSb \mid \epsilon$

- Derivation giving  $aaabbb$

## Examples

Grammar  $G_1$  with productions  $S \rightarrow aSb \mid \epsilon$

- Derivation giving  $aaabbb$



- Claim  $L(G_1) = \{a^n b^n \mid n \geq 0\}$  (How to prove?)

$\text{Prod}(G) S \xrightarrow{*} x \Rightarrow x \in a^m b^m$

Ind. on no of derivation steps:

( $n=1$ )  $S \xrightarrow{*} C$  clearly  $C \in L$

( $n=m$ )  $S \xrightarrow{*} a^m b^m \in L$

$a^n b \leftarrow L(c)$

Proof Ind. on n.

(n=0)  $C \rightarrow S$

(n=m+1)  $a^m b^m b$

$S \xrightarrow{m} a^m b^m$  Ind. Hyp

$aSb \xrightarrow{m} a a^m b^m b$

$S \xrightarrow{1} aSb \xrightarrow{m} a a^m b^m b$

## Examples

Grammar  $G_1$  with productions  $S \rightarrow aSb \mid \epsilon$

- Derivation giving  $aaabbb$
- Claim  $L(G_1) = \{ a^n b^n \mid n \geq 0 \}$  (How to prove?)

Grammar  $G_2$  with productions  $S \rightarrow aSa \mid bSb \mid a \mid b \mid \epsilon$

- Answer  $L(G_2) = \{ w \in \{a,b\}^* \mid w = x x^R \text{ or } x a x^R \text{ or } a x b x^R \}$

# Balanced Parenthesis Language

Example  $\boxed{\boxed{}}\boxed{\boxed{}}$  but not  $\boxed{\boxed{}}\boxed{\boxed{}}$ .

PAREN contains all strings  $x \in \{[], []\}^*$  s.t.

- ①  $\#[(x) = \#](x)$  count of left and right brackets is equal.
- ② for any prefix  $y$  of  $x$  we have  $\#[(y) \geq \#](y)$ .

$$\begin{array}{l} S \rightarrow S [S] S | C | SS \\ \searrow \\ [S] | SS | C \end{array}$$

# Balanced Parenthesis Language

Example  $\boxed{()}\boxed{)}$  but not  $\boxed{)}\boxed{()}$ .

PAREN contains all strings  $x \in \{[], []\}^*$  s.t.

- ①  $\#[(x)] = \#](x)$  count of left and right brackets is equal.
- ② for any prefix  $y$  of  $x$  we have  $\#[(y)] \geq \#](y)$ .

## Balanced Parenthesis Grammar

$$S \rightarrow SS \mid [S] \mid \epsilon$$

$$S \rightarrow SS \rightarrow [S]S$$

— —

Theorem  $PAREN = L(G)$

( $\exists$ )  $S \xrightarrow{*} x$  then  $x \in PAREN$   
 $S \xrightarrow{n} d$  then  $d$  satisfies (1) & (2),  $\forall n$ .  
proof by Ind. on length of derivations

Basis ( $n=0$ )  $S \xrightarrow{e} S$  check  $S$  satisfies (1) & (2)

Ind. Step:  $S \xrightarrow{n} \beta_1 S \beta_2 \xrightarrow{1} \alpha$

Ind Hyp:  $\beta_1 S \beta_2$  satisfies (1) & (2)

$\Rightarrow \beta_1 \beta_2$  also satisfies (1) & (2).

Case 1 ( $S \xrightarrow{e} e$ ):  $\therefore \alpha = \beta_1 \beta_2$  ✓

Case 2 ( $S \xrightarrow{SS}$ ):  $\therefore \alpha = \beta_1 S S \beta_2$

Case 3 ( $S \xrightarrow{[S]}$ ):  $\therefore \alpha = \beta_1 [S] \beta_2$

$\exists x \in \text{PAREN}$  then  $S \xrightarrow{*} x$ .

$[ [ ] ] | [ ]$

Prop: Ind on  $|x|$

Base ( $|x|=0$ )  $S \xrightarrow{*} G \checkmark$

y z

Ind. Step:  $x \in \text{PAREN}$

( $\Leftarrow$  1. If  $y$  proper prefix of  $x$ .  $|y| > 0$  s.t

$y$  satisfies (1) & (2)

i.e  $y = y \{ \dots \} \&$

$y \& z$  satisfies (1) & (2)

By Ind. Hyp  $S \xrightarrow{*} y \& S \xrightarrow{*} z$

$\therefore S \xrightarrow{*} S \xrightarrow{*} yS \xrightarrow{*} yz = x$

Case 2.  $x = [z]$

$\therefore z \in \text{PAREN}$

(Ind. Hyp.)  $S \xrightarrow{*} z$ .

$S \xrightarrow{*} [S] \xrightarrow{*} [z]$

- Defining **syntax** of programming languages.
- **Parsing** a given word.

Syntax directed processing.

See HMU Section 5.3 (Self Study).

# Backus Naur form grammar for Statement syntax in PL

```
<stmt> ::= <if-stmt> | <while-stmt> | <begin-stmt> | <assg-stmt>
<if-stmt> ::= if <bool-expr> then <stmt> else <stmt>
<while-stmt> ::= while <bool-expr> do <stmt>
<begin-stmt> ::= begin <stmt-list> end
<stmt-list> ::= <stmt> | <stmt> ; <stmt-list>
<assg-stmt> ::= <var> := <arith-expr>
<bool-expr> ::= <arith-expr><compare-op><arith-expr>
<compare-op> ::= < | > | ≤ | ≥ | = | ≠
<arith-expr> ::= <var> | <const> | (<arith-expr><arith-op><arith-expr>)
<arith-op> ::= + | - | * | /
<const> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
<var> ::= a | b | c | ⋯ | x | y | z
```

# Backus ~~Norm~~ Form Grammar for Statement Syntax in PL

*Naive*

$\langle \text{stmt} \rangle \xrightarrow{\quad}$

$\langle \text{while-stmt} \rangle \xrightarrow{\quad}$

$\text{while } \langle \text{bool-expr} \rangle \text{ do } \langle \text{stmt} \rangle \xrightarrow{\quad}$

$\text{while } \langle \text{arith-expr} \rangle \langle \text{compare-op} \rangle \langle \text{arith-expr} \rangle \text{ do } \langle \text{stmt} \rangle \xrightarrow{\quad}$

$\text{while } \langle \text{var} \rangle \langle \text{compare-op} \rangle \langle \text{arith-expr} \rangle \text{ do } \langle \text{stmt} \rangle$

$\text{while } \langle \text{var} \rangle \leq \langle \text{arith-expr} \rangle \text{ do } \langle \text{stmt} \rangle$

$\text{while } \langle \text{var} \rangle \leq \langle \text{var} \rangle \text{ do } \langle \text{stmt} \rangle$

$\text{while } x \leq \langle \text{var} \rangle \text{ do } \langle \text{stmt} \rangle$

$\text{while } x \leq y \text{ do } \langle \text{stmt} \rangle$

$\text{while } x \leq y \text{ do } \langle \text{begin-stmt} \rangle \xrightarrow{\quad}$

$\text{while } x \leq y \text{ do begin } x := \text{oct1}; y := y - 1$   
 $\text{end}$

Rightmost & Leftmost Derivation

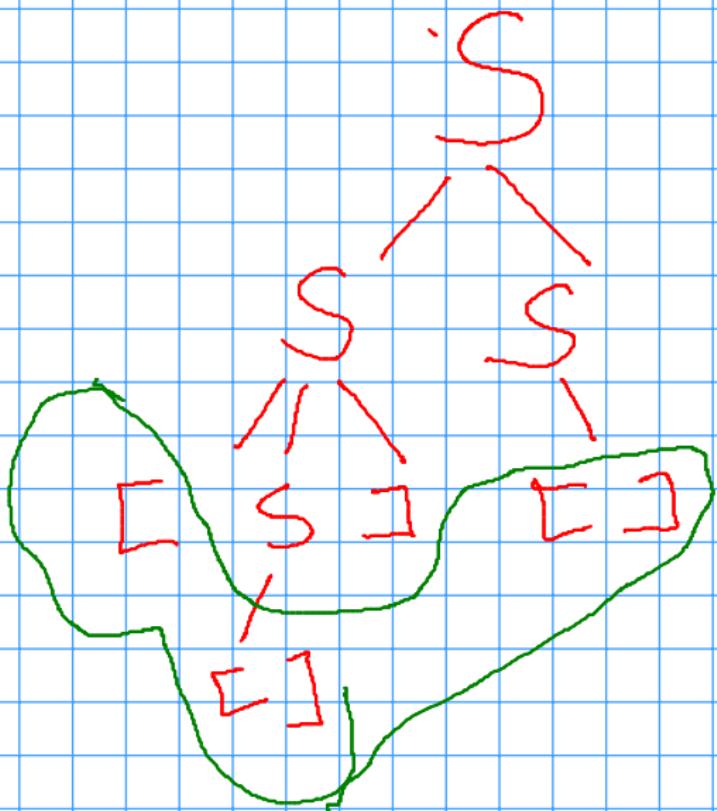


•  $S \rightarrow SS \rightarrow [S]S \rightarrow [ ]S \rightarrow [ ]\square$   
Leftmost Derivation

•  $S \rightarrow SS \rightarrow S[ ] \rightarrow [S][ ] \rightarrow [ ][ ]$   
Rightmost Derivation

•  $S \rightarrow SS \rightarrow [S]S \rightarrow [S][ ] \rightarrow [ ]\square$

# Parse Tree



# Context-Free Property.

- IF  $y_1 y_2 \dots y_k \xrightarrow{x} xc$  then  
 $\exists x_1, x_k. xc = x_1 \dots x_k$  and  
 $\forall i. y_i \rightarrow x_i$

• If  $\forall i, \gamma_i \xrightarrow{*} x_i$  then

$\gamma_1 \gamma_2 \cdots \gamma_k \xrightarrow{*} x_1 x_2 \cdots x_k$

$\gamma_1 \gamma_k \xrightarrow{*} x_1 \gamma_2 \gamma_k \xrightarrow{*} x_1 x_2 \gamma_3 \cdots \gamma_k$   
 $\xrightarrow{*} x_1 x_2 \cdots x_k$

SS

$S \xrightarrow{*} [[]]$

$S \xrightarrow{*} []$

$SS \rightarrow [[]][]$

# Transforming Given Grammar

- Preserve language :

Example Useless symbol A }

$$S \xrightarrow{*} \alpha A \beta.$$

- Remove all  $A \rightarrow \alpha$  rules

Useless 2  $S \xrightarrow{*} \alpha A \beta \xrightarrow{*} \alpha \gamma$

# Normal Forms for CFG

A CFG is in *Chomsky normal form* (CNF) if all productions are of the form

$$A \rightarrow BC \text{ or } A \rightarrow a,$$

where  $A, B, C \in N$  and  $a \in \Sigma$ . A CFG is in *Greibach normal form* (GNF) if all productions are of the form

$$A \rightarrow aB_1B_2 \cdots B_k$$

for some  $k \geq 0$ , where  $A, B_1, \dots, B_k \in N$  and  $a \in \Sigma$ . Note that  $k = 0$  is allowed, giving productions of the form  $A \rightarrow a$ .  $\square$

For example, the two grammars

$$S \rightarrow AB \mid AC \mid SS, \quad C \rightarrow SB, \quad A \rightarrow [, \quad B \rightarrow ], \quad (21.1)$$

$$S \rightarrow [B \mid [SB \mid [BS \mid [SBS, \quad B \rightarrow ] \quad (21.2)$$

## Theorem

For any CFG  $G$  we can construct  $G'$  in CNF and  $G''$  in GNF such that

$$L(G') = L(G'') = L(G) - \{\epsilon\}.$$

## Theorem

For any CFG  $G$  we can construct  $G'$  in CNF and  $G''$  in GNF such that

$$L(G') = L(G'') = L(G) - \{\epsilon\}.$$

- Epsilon production  $A \rightarrow \epsilon$  and unit production  $A \rightarrow B$ .

## Theorem

For any CFG  $G$  we can construct  $G'$  in CNF and  $G''$  in GNF such that

$$L(G') = L(G'') = L(G) - \{\epsilon\}.$$

- Epsilon production  $A \rightarrow \epsilon$  and unit production  $A \rightarrow B$ .
- Given  $G$  we construct  $G_1$  without epsilon or unit productions with  $L(G_1) = L(G) - \{\epsilon\}$

## Theorem

For any CFG  $G$  we can construct  $G'$  in CNF and  $G''$  in GNF such that

$$L(G') = L(G'') = L(G) - \{\epsilon\}.$$

- Epsilon production  $A \rightarrow \epsilon$  and unit production  $A \rightarrow B$ .
- Given  $G$  we construct  $G_1$  without epsilon or unit productions with  $L(G_1) = L(G) - \{\epsilon\}$
- We convert  $G_1$  into language equivalent  $G'$  in CNF.

## Theorem

For any CFG  $G$  we can construct  $G'$  in CNF and  $G''$  in GNF such that

$$L(G') = L(G'') = L(G) - \{\epsilon\}.$$

- Epsilon production  $A \rightarrow \epsilon$  and unit production  $A \rightarrow B$ .
- Given  $G$  we construct  $G_1$  without epsilon or unit productions with  $L(G_1) = L(G) - \{\epsilon\}$
- We convert  $G_1$  into language equivalent  $G'$  in CNF.
- We convert  $G'$  in CNF into language equivalent  $G''$  in GNF.

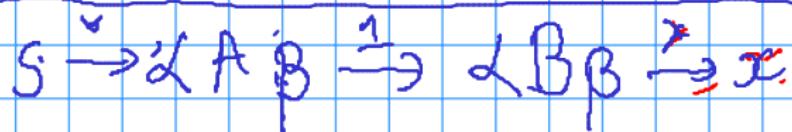
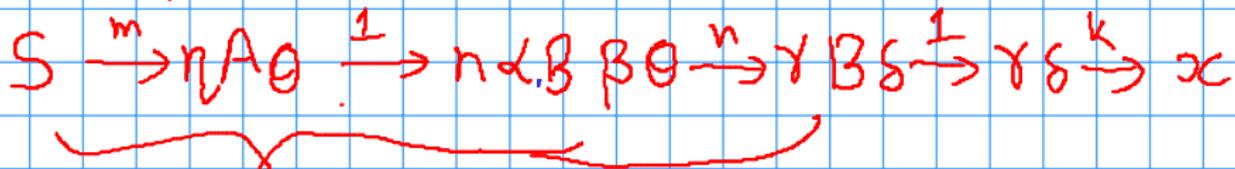
# Getting rid of epsilon and Unit Productions

Given  $G = (N, \Sigma, P, S)$  Construct  $\hat{P}$  by adding the following productions.

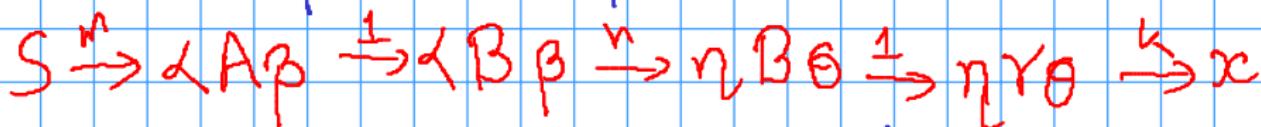
- (a) if  $A \rightarrow \alpha B \beta$  and  $B \rightarrow \epsilon$  are in  $\hat{P}$ , then  $A \rightarrow \alpha \beta$  is in  $\hat{P}$ ; and
  - (b) if  $A \rightarrow B$  and  $B \rightarrow \gamma$  are in  $\hat{P}$ , then  $A \rightarrow \gamma$  is in  $\hat{P}$ .
- 
- Let  $\hat{G} = (N, \Sigma, \hat{P}, S)$ .
  - **Claim** Let  $S \xrightarrow{*_{\hat{G}}} x$  via **shortest derivation**, where  $x \neq \epsilon$ .  
Then the derivation does not use  $\epsilon$  or unit productions.
  - Remove  $\epsilon$  and unit (and useless) productions from  $\hat{P}$  to get required grammar  $G_1$ .



$B \rightarrow \epsilon$



$A \rightarrow B$



$B \rightarrow \gamma$

# Example

Grammar  $S \rightarrow aSa \mid bSb \mid T, \quad T \rightarrow a \mid b \mid \epsilon.$

# Converting to CNF: An Example

Let  $G$  be  $S \rightarrow aSb \mid \epsilon$ .

# Converting to CNF: An Example

Let  $G$  be  $S \rightarrow aSb \mid \epsilon$ .

- We get  $G_1$  as  $S \rightarrow aSb \mid ab$  by removing epsilon and unit productions.

# Converting to CNF: An Example

Let  $G$  be  $S \rightarrow aSb \mid \epsilon$ .

- We get  $G_1$  as  $S \rightarrow aSb \mid ab$  by removing epsilon and unit productions.
- Replace terminals by fresh nonterminals.

$S \rightarrow ASB \mid AB, \quad A \rightarrow a, \quad B \rightarrow b.$

# Converting to CNF: An Example

Let  $G$  be  $S \rightarrow aSb \mid \epsilon$ .

- We get  $G_1$  as  $S \rightarrow aSb \mid ab$  by removing epsilon and unit productions.

- Replace terminals by fresh nonterminals.

$$S \rightarrow ASB \mid AB, \quad A \rightarrow a, \quad B \rightarrow b.$$

- Replace long RHS to get  $G'$  in CNF.

$$S \rightarrow AC \mid AB, \quad C \rightarrow SB, \quad A \rightarrow a, \quad B \rightarrow b$$

# Converting to CNF: An Example

Let  $G$  be  $S \rightarrow aSb \mid \epsilon$ .

- We get  $G_1$  as  $S \rightarrow aSb \mid ab$  by removing epsilon and unit productions.
- Replace terminals by fresh nonterminals.

$S \rightarrow ASB \mid AB, \quad A \rightarrow a, \quad B \rightarrow b.$

- Replace long RHS to get  $G'$  in CNF.

$S \rightarrow AC \mid AB, \quad C \rightarrow SB, \quad A \rightarrow a, \quad B \rightarrow b$

It is easy to see that  $L(G') = L(G_1)$ .

# Converting CNF to GNF

Omitted.

# Some properties of CNF and GNF

- Every string in the language can be derived by left-most (or right-most) derivation.
  - If  $S \xrightarrow{G}^n x$  and  $|x| = m$  then
    - If  $G$  is CNF we have  $n = ?$  (in terms of  $m$ )
    - If  $G$  is GNF we have  $n = ?$  (in terms of  $m$ )
- non-term + 1  
each non-term  
= 1 term  
 $\xrightarrow{n}$   
 $n = 2m - 1$*

*for GNF  $\rightarrow n = m$  (each step, L++)*

# CFL and non-CFL languages

We gave CFG for following languages to show that they are CFL

- $a^n b^n$
- PALINDROME  $\{x \mid x = \text{rev}(x)\}$
- Balanced Paranthesis language PAREN

# CFL and non-CFL languages

We gave CFG for following languages to show that they are CFL

- $a^n b^n$
- PALINDROME  $\{x \mid x = \text{rev}(x)\}$
- Balanced Paranthesis language PAREN
- What about  $a^n b^n c^n$ ?

# CFL and non-CFL languages

We gave CFG for following languages to show that they are CFL

- $a^n b^n$
- PALINDROME  $\{x \mid x = \text{rev}(x)\}$
- Balanced Paranthesis language PAREN
- What about  $a^n b^n c^n$ ?
- This language is not context free. How to prove this?

# CFL and non-CFL languages

We gave CFG for following languages to show that they are CFL

- $a^n b^n$
- PALINDROME  $\{x \mid x = \text{rev}(x)\}$
- Balanced Paranthesis language PAREN
- What about  $a^n b^n c^n$ ?
- This language is not context free. How to prove this?
- Answer: Pumping Lemma for Context Free Languages.

- Suppose  $L$  is CFL. Hence a CNF grammar  $G$  generates it.

- Suppose  $L$  is CFL. Hence a CNF grammar  $G$  generates it.
- Consider a **parse tree**  $T_x$  for a **word**  $x$  obtained using  $G$ .

- Suppose  $L$  is CFL. Hence a CNF grammar  $G$  generates it.
- Consider a **parse tree**  $T_x$  for a **word**  $x$  obtained using  $G$ .
- If  $x$  is long (i.e.  $|x| \geq 2^{n+1}$ ) then  $T_x$  has depth at least  $n + 1$ .

- Suppose  $L$  is CFL. Hence a CNF grammar  $G$  generates it.
- Consider a **parse tree**  $T_x$  for a **word**  $x$  obtained using  $G$ .
- If  $x$  is long (i.e.  $|x| \geq 2^{n+1}$ ) then  $T_x$  has depth at least  $n + 1$ .
- There is a long path with at least  $n + 1$  internal nodes from root to leaf.

- Suppose  $L$  is CFL. Hence a CNF grammar  $G$  generates it.
- Consider a **parse tree**  $T_x$  for a **word**  $x$  obtained using  $G$ .
- If  $x$  is long (i.e.  $|x| \geq 2^{n+1}$ ) then  $T_x$  has depth at least  $n + 1$ .
- There is a long path with at least  $n + 1$  internal nodes from root to leaf.
- Let  $n$  be the number of NONTERMINALS in  $G$ .

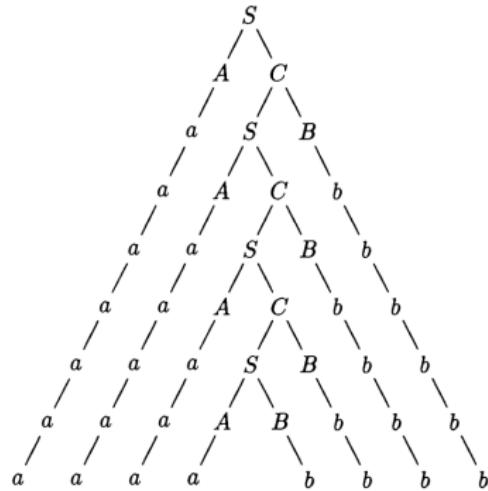
- Suppose  $L$  is CFL. Hence a CNF grammar  $G$  generates it.
- Consider a **parse tree**  $T_x$  for a **word**  $x$  obtained using  $G$ .
- If  $x$  is long (i.e.  $|x| \geq 2^{n+1}$ ) then  $T_x$  has depth at least  $n + 1$ .
- There is a long path with at least  $n + 1$  internal nodes from root to leaf.
- Let  $n$  be the number of NONTERMINALS in  $G$ .
- Hence some non-terminal repeats on the long path.

- Suppose  $L$  is CFL. Hence a CNF grammar  $G$  generates it.
- Consider a **parse tree**  $T_x$  for a **word**  $x$  obtained using  $G$ .
- If  $x$  is long (i.e.  $|x| \geq 2^{n+1}$ ) then  $T_x$  has depth at least  $n + 1$ .
- There is a long path with at least  $n + 1$  internal nodes from root to leaf.
- Let  $n$  be the number of NONTERMINALS in  $G$ .
- Hence some non-terminal repeats on the long path.
- Consider lowest (closest to leaf) pair of repeating non-terminals on the long path.

- Suppose  $L$  is CFL. Hence a CNF grammar  $G$  generates it.
- Consider a **parse tree**  $T_x$  for a **word**  $x$  obtained using  $G$ .
- If  $x$  is long (i.e.  $|x| \geq 2^{n+1}$ ) then  $T_x$  has depth at least  $n + 1$ .
- There is a long path with at least  $n + 1$  internal nodes from root to leaf.
- Let  $n$  be the number of NONTERMINALS in  $G$ .
- Hence some non-terminal repeats on the long path.
- Consider lowest (closest to leaf) pair of repeating non-terminals on the long path.
- This can be pumped arbitrary  $i$  times.

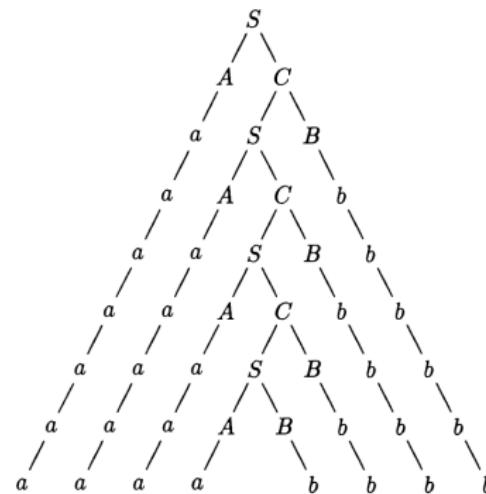
# Parse Tree of CNF Grammer Derivation

Let  $G$  be  $S \rightarrow AC \mid AB, \quad C \rightarrow SB, \quad A \rightarrow a, \quad B \rightarrow b.$



# Parse Tree of CNF Grammer Derivation

Let  $G$  be  $S \rightarrow AC \mid AB, \quad C \rightarrow SB, \quad A \rightarrow a, \quad B \rightarrow b.$



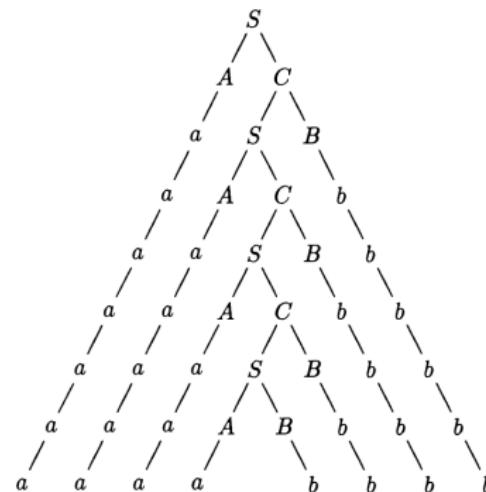
## Leftmost Derivation

$$\begin{aligned} S &\rightarrow AC \rightarrow aC \rightarrow aSB \rightarrow aACB \rightarrow aaCB \rightarrow aaSBB \rightarrow aaACBB \\ &\rightarrow aaaCBB \rightarrow aaaSBBB \rightarrow aaaABB BB \rightarrow aaaaBBBB \\ &\rightarrow aaaabBBB \rightarrow aaaabbBB \rightarrow aaaabbbB \rightarrow aaaabbbb, \end{aligned}$$



# Parse Tree of CNF Grammer Derivation

Let  $G$  be  $S \rightarrow AC \mid AB, \quad C \rightarrow SB, \quad A \rightarrow a, \quad B \rightarrow b.$



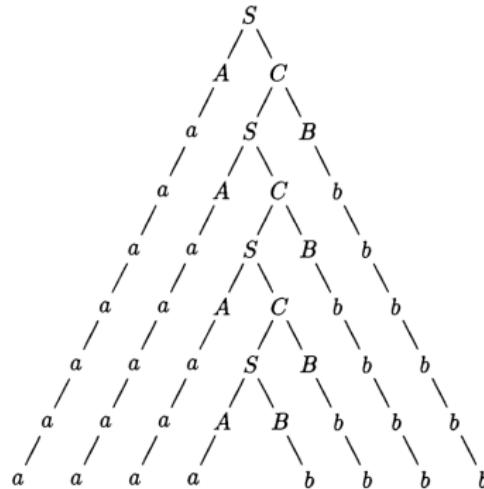
## Rightmost Derivation

$$\begin{aligned} S &\rightarrow AC \rightarrow ASB \rightarrow ASb \rightarrow AACb \rightarrow AASBb \rightarrow AASbb \rightarrow AAACbb \\ &\rightarrow AAASBbb \rightarrow AAASbbb \rightarrow AAAABbbb \rightarrow AAAAbbbb \\ &\rightarrow AAAabbbb \rightarrow AAaabbbb \rightarrow Aaaabbbb \rightarrow aaaabbbb, \end{aligned}$$

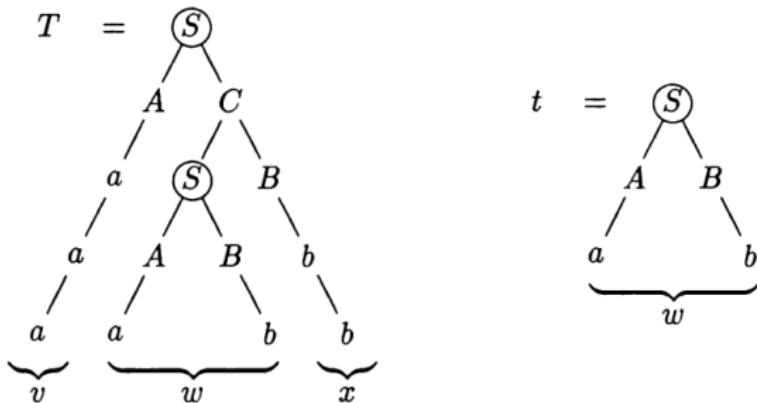


# Parse Tree of CNF Grammer Derivation

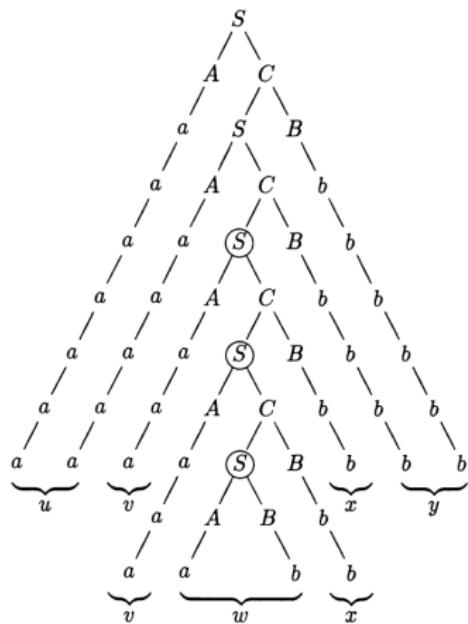
Let  $G$  be  $S \rightarrow AC \mid AB, \quad C \rightarrow SB, \quad A \rightarrow a, \quad B \rightarrow b.$



# Subtrees for lowermost pair of repeating Nonterminals



Subtrees for lowermost pair of repeating Nonterminals



# Pumping Lemma for CFL

**Lemma** If  $A$  is CFL then

$\exists k > 0$  such that

$\forall z \in A$  with  $|z| \geq k$ ,

$\exists u, v, w, x, Y$  with  $z = u \cdot v \cdot w \cdot x \cdot y$  and

$|vwx| \leq k$  and  $|vx| > 0$  s.t.

$u \cdot (v^i) \cdot w \cdot (x^j) \cdot y \in A$  for all  $i \geq 0$ .

# Pumping Lemma in Contrapositive Form

Game with demon to prove that Language  $A$  is not CFL

- Demon chooses arbitrary  $k > 0$ .
- You choose  $z \in A$  with  $|z| \geq k$ ,
- Demon chooses  $u, v, w, x, Y$  with  $z = u \cdot v \cdot w \cdot x \cdot y$  and  $|vwx| \leq k$  and  $|vx| > 0$
- You choose  $i \geq 0$  and show that  $u \cdot (v^i) \cdot w \cdot (x^i) \cdot y \notin A$  to win the game.

You give a winning strategy for the above game to show that  $A$  is not CFL.

## example

$a^n b^n c^n$  is not CFL. Proof using pumping lemma.

- Demon chooses  $k > 0$ .
- Choose  $z = a^k b^k c^k$  which is in  $A$ .
- Demon splits  $z = u \cdot v \cdot w \cdot x \cdot y$  with conditions. Some Cases:
  - ①  $u = \epsilon$  and  $vwx = a^k$  and  $y = b^k c^k$ .  
You choose  $i = 2$  which given  $v^2wx^2 = a^j$  with  $j > k$ . Hence,  
 $u \cdot (v^2) \cdot w \cdot (x^2) \cdot y \notin A$ .
  - ② Demon chooses  $vwx$  within  $a^k$  part. Choosing  $i = 2$  and  
pumping up, increases count of  $a$  without changing counts of  $b$   
and  $c$ . Hence pumped word is not in  $A$ .
  - ③ Demon chooses  $vwx = a^j b^l$ . You choose  $i = 2$ . Pumping up  
increases the count of  $a +$  count of  $b$  without changing the  
count of  $c$ . Hence the pumped word is not in  $A$ .
  - ④ The Remaining cases are similar.

## example

Let  $A = \{w \cdot w \mid w \in \{a, b\}^*\}$ . Prove that  $A$  is not CFL.

## Theorem

If  $A, B$  are any CFL languages then

- $A \cup B$  is CFL
  - $\text{rev}(A)$  is CFL
  - $A \cdot B$  is CFL
  - $A^*$  is CFL
- $h$  is homomorphism from  $\Sigma^* \rightarrow \Delta^*$   
then  $h(A)$  is CFL
- $h^{-1}(A)$  is CFL

# Closure Property Constructors

Let  $G_1 = (N_1, \Sigma, P_1, S_1)$  and  $G_2 = (N_2, \Sigma, , P_2, S_2)$   
with  $N_1 \cap N_2 = \emptyset$ .

# Closure Property Constructions

Let  $G_1 = (N_1, \Sigma, P_1, S_1)$  and  $G_2 = (N_2, \Sigma, P_2, S_2)$   
with  $N_1 \cap N_2 = \emptyset$ .

- Let  $G = (N_1 \cup N_2 \cup \{S\}, \Sigma, P, S)$  with  
 $P = \{S \rightarrow S_1 \mid S_2\} \cup P_1 \cup P_2$ .  
Then,  $L(G) = L(G_1) \cup L(G_2)$ . How?

# Closure Property Constructions

Let  $G_1 = (N_1, \Sigma, P_1, S_1)$  and  $G_2 = (N_2, \Sigma, P_2, S_2)$   
with  $N_1 \cap N_2 = \emptyset$ .

- Let  $G = (N_1 \cup N_2 \cup \{S\}, \Sigma, P, S)$  with  
 $P = \{S \rightarrow S_1 | S_2\} \cup P_1 \cup P_2$ .  
Then,  $L(G) = L(G_1) \cup L(G_2)$ . How?
- Let  $G = (N_1 \cup N_2 \cup \{S\}, \Sigma, P, S)$  with  
 $P = \{S \rightarrow S_1 S_2\} \cup P_1 \cup P_2$ .  
Then,  $L(G) = L(G_1) \cdot L(G_2)$ .
- $L(G) = (L(G_1))^*$

$$S \rightarrow S_1 | S_2,$$

$$\left. \begin{array}{l} L(G) = L(G_1) \\ A \xrightarrow{d} C - G_1 \\ \text{and} \\ A \xrightarrow{d} L \end{array} \right\} R$$

$$S \rightarrow aSb \mid C$$
$$a^n b^n$$
$$S \rightarrow bSa \mid C$$
$$b^n a^n$$

# Closure Property Constructions

Let  $G_1 = (N_1, \Sigma, P_1, S_1)$  and  $G_2 = (N_2, \Sigma, P_2, S_2)$   
with  $N_1 \cap N_2 = \emptyset$ .

- Let  $G = (N_1 \cup N_2 \cup \{S\}, \Sigma, P, S)$  with  
 $P = \{S \rightarrow S_1 | S_2\} \cup P_1 \cup P_2$ .  
Then,  $L(G) = L(G_1) \cup L(G_2)$ . How?
- Let  $G = (N_1 \cup N_2 \cup \{S\}, \Sigma, P, S)$  with  
 $P = \{S \rightarrow S_1 S_2\} \cup P_1 \cup P_2$ .  
Then,  $L(G) = L(G_1) \cdot L(G_2)$ .
- Let  $G = (N_1 \cup \{S\}, \Sigma, P, S)$  with  $P = \{S \rightarrow S_1 S | \epsilon\} \cup P_1$ .  
Then,  $L(G) = (L(G_1))^*$ .

CFL are not closed under  $\cap$

$a^nb^n$

$a^nb^n$

If  $A, B$  are CFL languages then Intersection  $A \cap B$  need not be CFL.

$a^n b^n c^m$

$a^m b^n c^n$

$\{a^nb^n | n \geq 0\}$

## CFL are not closed under $\cap$

If  $A, B$  are CFL languages then Intersection  $A \cap B$  need not be CFL.

- How to show non-closure?

# CFL are not closed under $\cap$

If  $A, B$  are CFL languages then Intersection  $A \cap B$  need not be CFL.

- How to show non-closure?
- Give a **counter example**. Choose CFL  $A, B$  s.t.  $A \cap B$  is not CSL.

## CFL are not closed under $\cap$

If  $A, B$  are CFL languages then Intersection  $A \cap B$  need not be CFL.

- How to show non-closure?
- Give a **counter example**. Choose CFL  $A, B$  s.t.  $A \cap B$  is not CSL.

Let  $A = \{a^n b^n c^m \mid m, n \geq 0\}$ , and  $B = \{a^m b^n c^n \mid m, n \geq 0\}$

## CFL are not closed under $\cap$

If  $A, B$  are CFL languages then Intersection  $A \cap B$  need not be CFL.

- How to show non-closure?
- Give a **counter example**. Choose CFL  $A, B$  s.t.  $A \cap B$  is not CSL.

Let  $A = \{a^n b^n c^m \mid m, n \geq 0\}$ , and  $B = \{a^m b^n c^n \mid m, n \geq 0\}$

- Show that  $A$  as well as  $B$  are CFL.

## CFL are not closed under $\cap$

If  $A, B$  are CFL languages then Intersection  $A \cap B$  need not be CFL.

- How to show non-closure?
- Give a **counter example**. Choose CFL  $A, B$  s.t.  $A \cap B$  is not CSL.

Let  $A = \{a^n b^n c^m \mid m, n \geq 0\}$ , and  $B = \{a^m b^n c^n \mid m, n \geq 0\}$

- Show that  $A$  as well as  $B$  are CFL.
- $A \cap B = \{a^n b^n c^n \mid n \geq 0\}$ . Show that this is not CFL.

Proof by contradiction

- Assume to contrary that for every CFL  $A$  the language  $\sim A$  is also CFL.

Proof by contradiction

- Assume to contrary that for every CFL  $A$  the language  $\sim A$  is also CFL.
- Let  $A, B$  be arbitrary CFL.

## Proof by contradiction

- Assume to contrary that for every CFL  $A$  the language  $\sim A$  is also CFL.
- Let  $A, B$  be arbitrary CFL.
- By set theory  $A \cap B = \sim ((\sim A) \cup (\sim B))$ .

## Proof by contradiction

- Assume to contrary that for every CFL  $A$  the language  $\sim A$  is also CFL.
- Let  $A, B$  be arbitrary CFL.
- By set theory  $A \cap B = \sim ((\sim A) \cup (\sim B))$ .
- Using closure under  $\cup$  and  $\sim$ , then RHS is also CFL.

## Proof by contradiction

- Assume to contrary that for every CFL  $A$  the language  $\sim A$  is also CFL.
- Let  $A, B$  be arbitrary CFL.
- By set theory  $A \cap B = \sim ((\sim A) \cup (\sim B))$ .
- Using closure under  $\cup$  and  $\sim$ , then RHS is also CFL.
- This contradicts that CFL is not closed under  $\cap$ .

# Decision Problems

Given CFG  $G$  and  $x \in \Sigma^*$ ,

- We can check by an algorithm whether  $x \in L(G)$ .  
CYK algorithm (omitted).

Given CFG  $G$  and  $x \in \Sigma^*$ ,

- We can check by an algorithm whether  $x \in L(G)$ .  
CYK algorithm (omitted).
- **Parsing:** Constructing parse tree for given word  $x$  under CFG  $G$ .

Widely used in compiler design.

Specialized grammars such as LALR, SLR or LR( $k$ ) have efficient parsing algorithms. Aho, Ullman: Turing Award 2020.

## Decision Problems

Given CFG  $G$  and  $x \in \Sigma^*$ ,

- We can check by an algorithm whether  $x \in L(G)$ .  
CYK algorithm (omitted).
  - **Parsing**: Constructing parse tree for given word  $x$  under CFG  $G$ .

Widely used in compiler design.

Specialized grammars such as LALR, SLR or LR( $k$ ) have efficient parsing algorithms. Aho, Ullman: Turing Award 2020.

- We can check by an algorithm whether  $L(G) = \emptyset$ .  
Emptiness of CFG is decidable.

Given  $G_1$  and  $G_2$

$L(A_1) \subseteq L(A_2)$ ?

## Inclusion

# Decision Problems

Given CFG  $G$  and  $x \in \Sigma^*$ ,

- We can check by an algorithm whether  $x \in L(G)$ .  
CYK algorithm (omitted).
- **Parsing:** Constructing parse tree for given word  $x$  under CFG  $G$ .  
Widely used in compiler design.  
Specialized grammars such as LALR, SLR or LR( $k$ ) have efficient parsing algorithms. Aho, Ullman: Turing Award 2020.
- We can check by an algorithm whether  $L(G) = \emptyset$ .  
**Emptiness of CFG is decidable.**
- There is no algorithm to decide whether  $L(G) = \Sigma^*$ .  
**Universality of CFG is undecidable.**

