

# INT201 Decision, Computation and Language

Lecture 6 – Context-Free Languages (1)

Dr Yushi Li



# Context-Free Languages

- Context-Free Grammar (CFG)
- Chomsky Normal Form (CNF)



# Context-Free Languages

自动机

接受语言的字符串

- **Finite automata accept** precisely the strings in the language.

通过计算去判断输入的语言是否属于该语言

*Perform a computation to determine whether a specific string is in the language.*

正规化表达

描述语言的字符串

- **Regular expressions describe** precisely the strings in the language

*Describe the general shape of all strings in the language.*

classify

- **Context-free grammar (CFG)** is an entirely different formalism for defining a class of languages.

定义语言的种类

*Give a procedure for listing off all strings in the language.*

提供一个步骤来列出语言中所有strings



# Context-Free Languages 上下文无关语言

## Applications of CFG

- Programming languages: CFGs are used to define the syntax of programming languages, allowing parsers to analyze code structure.  
自然语言处理 分析程序 分析代码框架
- NLP: CFGs help in parsing sentences, enabling applications like machine translation and speech recognition.  
编译器 解析句子
- Compilers: CFGs facilitate syntax analysis, ensuring that the source code adheres to the language's grammatical rules.  
确保源代码 贴合 语言 语法规则



# Context-Free Grammar

## Example

- Start variable S with rules:

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$A \rightarrow aA$$

$$B \rightarrow b$$

$$B \rightarrow bB$$

$$L = \{a^m | b^m : m \geq 1\}$$

variables: S, A, B terminals: a, b

- Following these rules, we can yield ?

We can infer a language from given rule

$$S \Rightarrow AB \Rightarrow aAB \Rightarrow aAbB \Rightarrow aaAbB \Rightarrow$$

$$aaaAbB \Rightarrow aaaAbbB \Rightarrow aaaabb .$$

aaa  
bbb  
aabbb... .



# Context-Free Grammar

## Definition

A context-free grammar is a 4-tuple  $G = (V, \Sigma, R, S)$ , where

1.  $V$  is a finite set, whose elements are called **variables**,  
有向集合
2.  $\Sigma$  is a finite set, whose elements are called **terminals**,  
未端 (注意和 DFA / NFA 的区别)
3.  $V \cap \Sigma = \emptyset$ ,    variable  $\cap$  terminal    没有元素相交
4.  $S$  is an element of  $V$ ; it is called the **start variable**,  
开始
5.  $R$  is a finite set, whose elements are called **rules**. Each rule has the form  $A \rightarrow w$ ,  
where  $A \in V$  and  $w \in (V \cup \Sigma)^*$ .

$A$  is a variable in  $V$

$w$  is the strings constructed from  $(V \cup \Sigma)^*$



# Context-Free Grammar

## Example

$0^k | 1^k$

Language  $L = \{0^k | 1^k : k \geq 0\}$  has CFG  $G = (V, \Sigma, R, S)$ ,

variable set  $V = \{S\}$

Terminal set  $\Sigma = \{0, 1\}$

start variable  $S$

Rule set  $R: S \rightarrow 0S1$

$0^k | 1^k$

$S \rightarrow \epsilon$

$||$

$S \rightarrow 0(S)1 \rightarrow 00S11 \rightarrow 000S111 \Rightarrow 0\dots 0S1\dots 1$



# Deriving strings and languages using CFG

$\Rightarrow : \text{yeild}$

产出

Let  $G = (V, \Sigma, R, S)$  be a context free grammar with

- $A \in V$
- $u, v, w \in (V \cup \Sigma)^*$ ,
- $A \rightarrow w$  is a rule of the grammar

The string  $uvw$  can be derived in one step from the string  $uAv$ , written as

$$\underline{uAv} \Rightarrow \underline{uvw}$$

**Example:** aaAbb  $\Rightarrow$  aaaAbb



# Deriving strings and languages using CFG

$\xrightarrow{*}$  : derive

右由左得到

Let  $G = (V, \Sigma, R, S)$  be a context free grammar with

- $u, v \in (V \cup \Sigma)^*$

得到

The string v can be derived from the string u, written as  $u \xrightarrow{*} v$ , if one of the following conditions holds:

1.  $u = v$

2. there exist an integer  $k \geq 2$  and a sequence  $u_1, u_2, \dots, u_k$  of strings in  $(V \cup \Sigma)^*$ , such that

(a)  $u = u_1,$

$$u_1, u_2, \dots, u_k \in (V \cup \Sigma)^*$$

(b)  $v = u_k$ , and  $u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_k$ .

**Example:** With the rules  $A \rightarrow B1 \mid D0C$

$$0AA \xrightarrow{*} 0D0CB1$$



# Language of CFG

## Definition

The language of CFG  $G = (V, \Sigma, R, S)$  is

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

Such a language is called **context-free**, and satisfies  $L(G) \subseteq \Sigma^*$ .

## Example

CFG  $G = (V, \Sigma, R, S)$  with

1.  $V = \{S\}$

2.  $\Sigma = \{0, 1\}$

3. Rules  $R$ :  $S \rightarrow 0S \mid \epsilon$

$L(G) = ?$

$$\begin{aligned} & S \xrightarrow{\text{0S}} 0S \rightarrow 00S \rightarrow \dots \\ & \quad \rightarrow 0 \dots 0S \\ & \therefore S \xrightarrow{\epsilon} \end{aligned}$$

$$\therefore S \xrightarrow{\epsilon} 0 \dots 0 \Rightarrow L(G) = \{0^n : n \geq 0\}$$



## Example (Palindrome) 回文

CFG  $G = (V, \Sigma, R, S)$  with

1.  $V = \{S\}$
2.  $\Sigma = \{a, b\}$
3. Rules  $R: S \rightarrow aSa \mid bSb \mid a \mid b \mid \epsilon$

$$\left\{ \begin{array}{l} S \rightarrow aSa \\ S \rightarrow bSb \\ S \rightarrow a \\ S \rightarrow b \\ S \rightarrow \epsilon \end{array} \right.$$

Language of this CFG ?

$$S \Rightarrow aSa \Rightarrow aaSaa \xrightarrow{*} a \dots a S a \dots a$$

$$\Rightarrow \left\{ \begin{array}{ll} a \dots aaa \dots a & S \rightarrow a \\ a \dots aba \dots a & S \rightarrow b \\ a \dots a \epsilon a \dots a & S \rightarrow \epsilon \end{array} \right.$$

$$S \Rightarrow bSb \Rightarrow bbSbb \dots \text{ same measure as above}$$

$$L(G) = \{w \in \Sigma^* \mid w = w^R\} \quad R: \text{reverse}$$



# 简单的算术表达

## Example (Simple Arithmetic Expressions)

CFG  $G = (V, \Sigma, R, S)$  with

1.  $V = \{S\}$
2.  $\Sigma = \{+, -, \times, /, (, ), 0, 1, 2, \dots, 9\}$
3. Rules  $R$ :  
$$S \rightarrow S + S \mid S - S \mid S \times S \mid S / S \mid (S) \mid -S \mid 0 \mid 1 \mid \dots \mid 9$$

$L(G)$ : valid arithmetic expressions over single-digit integers

$S$  derives string  $3 \times (5 + 6)$ ?

$$\begin{aligned} S &\Rightarrow S \times S \Rightarrow S \times (S) \Rightarrow S \times (S + S) \Rightarrow 3 \times (S + S) \Rightarrow 3 \times (S + S) \\ &\Rightarrow 3 \times (5 + 6) \end{aligned}$$



# Regular Languages are context-free

(if)

(could say)

Theorem Regular Language  $\Rightarrow$  Context free

Let  $\Sigma$  be an alphabet and let  $L \subseteq \Sigma^*$  be a regular language. Then  $L$  is a context-free language (Every regular language is context-free).

Proof (general idea)

L是正则语言，有一个DFA M接受，L是上下文无关的需要  
有一个上下文无关语法 G 满足  $L = L(M) = L(G)$

Since  $L$  is a regular language, there exists a deterministic finite automaton  $M = (Q, \Sigma, \delta, q_0, F)$  that accepts  $L$ . To prove that  $L$  is context-free, we have to define a context-free grammar  $G = (V, \Sigma, R, S)$ , such that  $L = L(M) = L(G)$ . Thus,  $G$  must have the following property:

$$w \in L(M) \Leftrightarrow w \in L(G)$$

For every string  $w \in \Sigma^*$ ,

$$\underline{w} \in L(M) \text{ if and only if } w \in L(G),$$

which can be reformulated as

$$M \text{ accepts } w \text{ if and only if } \underline{S \xrightarrow{*} w}.$$

$G$ 的  $V$  就是从的所有的  $Q$

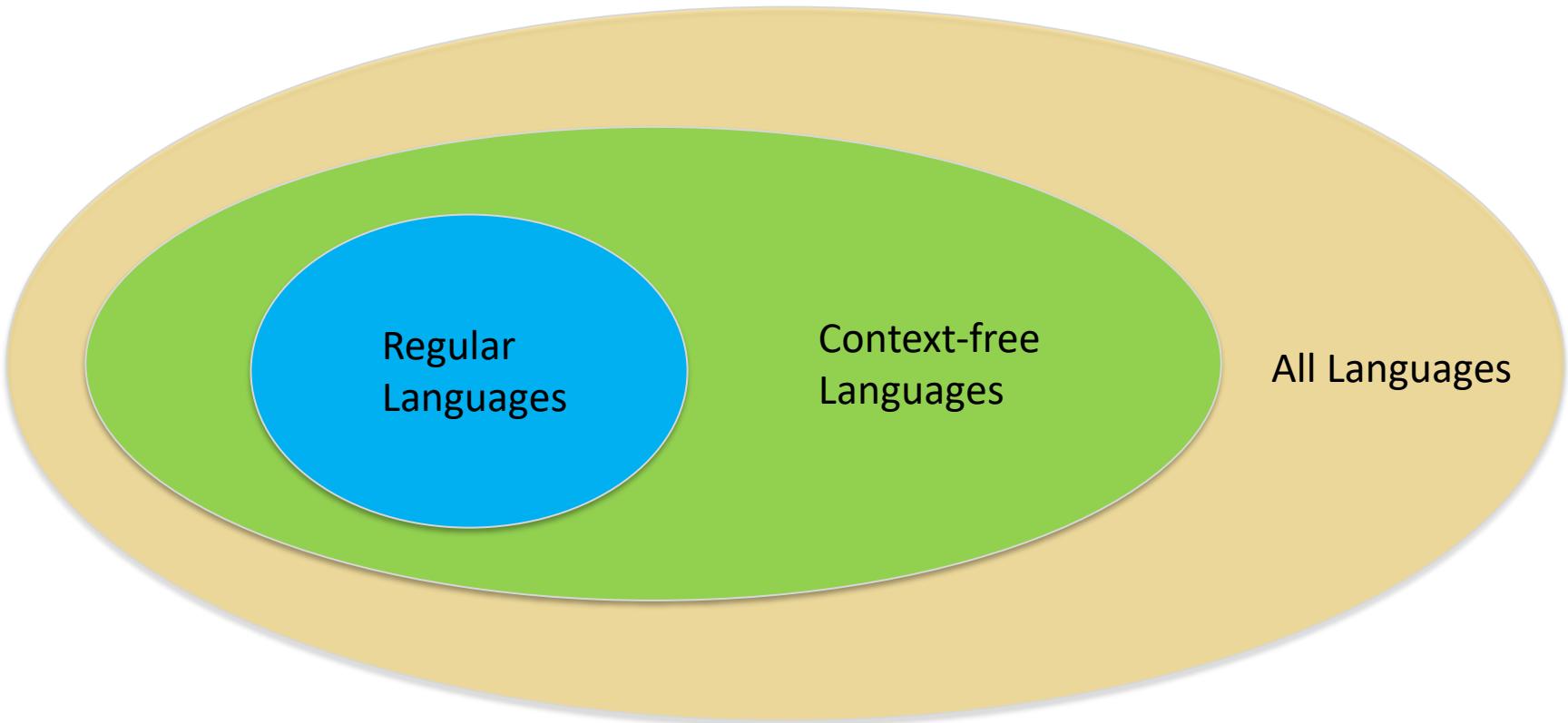
Set  $\underline{V} = \{R_i \mid q_i \in Q\}$  (that is,  $G$  has a variable for every state of  $M$ ). Now, for every transition  $\underline{\delta(q_i, a) = q_j}$  add a rule  $R_i \rightarrow aR_j$ . For every accepting state  $\underline{q_i \in F}$  add a rule  $\underline{R_i \rightarrow \varepsilon}$ . Finally, make the start variable  $S = R_0$ .



$R_0$  is the initial state of the machine

Regular Languages are context-free

$L \text{ is regular} \Rightarrow L \text{ is context free}$   
 $\Leftarrow$



**Closure properties of CFLs:** CFLs are closed under operations like union and concatenation but not under intersection or complementation.



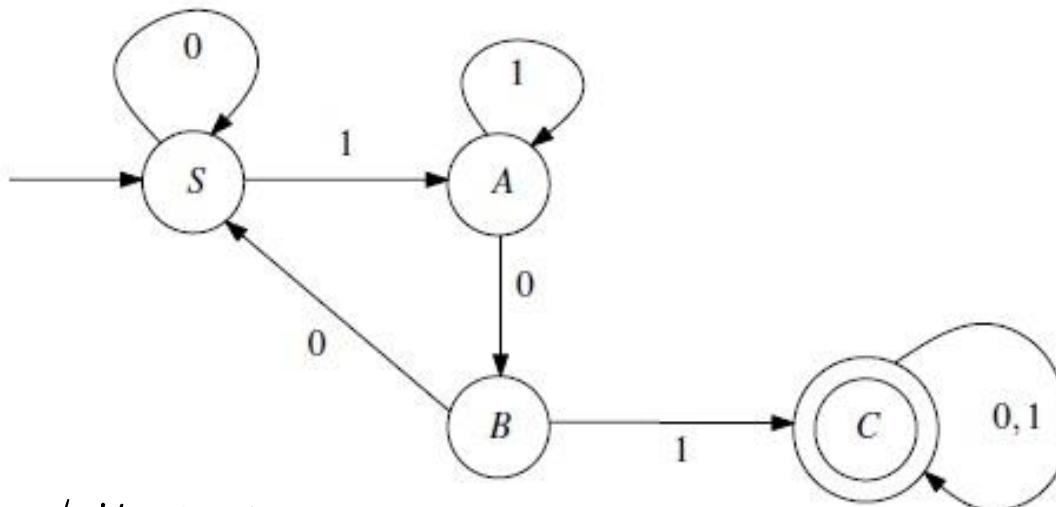
# Regular Languages are context-free

## Example

Let  $L$  be the language defined as

$$L = \{w \in \{0, 1\}^*: \underline{\text{101}} \text{ is a substring of } w\}.$$

The DFA  $M$  that accepts  $L$



将DFA转换为CFG

How can we convert  $M$  to a context-free grammar  $G$  whose  $language$  is  $L$ ?



# Regular Languages are context-free

## Example

- $G = \{V, \Sigma, R, S\}$

$$V = \{S, A, B, C\}$$

$$\Sigma = \{0, 1\}$$

Start variable :  $S$  (initial state of M)

Rules:

$$S \rightarrow 0S \mid 1A$$

$$A \rightarrow 0B \mid 1A$$

$$B \rightarrow 0S \mid 1C$$

$$C \rightarrow 0C \mid 1C \mid \epsilon$$



# Chomsky Normal Form (CNF) 克姆斯基公式

## Definition

A context-free grammar  $G = (V, \Sigma, R, S)$  is said to be in **Chomsky normal form**, if every rule in  $R$  has one of the following three forms: *如果 Rules 滿足下面三點條件*

- $A \rightarrow BC$ , where  $A, B$ , and  $C$  are elements of  $V$ ,  $B \neq S$ , and  $C \neq S$ .
- $A \rightarrow a$ , where  $A$  is an element of  $V$  and  $a$  is an element of  $\Sigma$ .
- $S \rightarrow \epsilon$ , where  $S$  is the start variable.

## Why CNF?

Grammars in Chomsky normal form are far easier to analyze.

## Example

Rules of CFG in Chomsky normal form with  $V = \{S, A, B\}$ ,  $\Sigma = \{a, b\}$ :

$G_1 : S \rightarrow AB, S \rightarrow c, A \rightarrow a, B \rightarrow b$  (CNF)

$G_1 : S \rightarrow aA, A \rightarrow a, B \rightarrow c$  (not CNF)



# Chomsky Normal Form (CNF)

## Theorem

Let  $\Sigma$  be an alphabet and let  $L \subseteq \Sigma^*$  be a context-free language. There exists a context-free grammar in Chomsky normal form, whose language is  $L$  (Every CFL can be described by a CFG in CNF).

---

## CFL $\rightarrow$ CNF

Given CFG  $G = (V, \Sigma, R, S)$ . Replace, one-by-one, every rule that is not “Chomsky”.

- Start variable (not allowed on RHS of rules)
- $\epsilon$ -rules ( $A \rightarrow \epsilon$  not allowed when  $A$  isn't start variable)
- all other violating rules ( $A \rightarrow B, A \rightarrow aBc, A \rightarrow BCDE$ )



context free grammar  $\rightarrow$  chomsky normal form

## Converting CFG into CNF

### Transformation steps

1. ~~从右端除掉 start variable~~

Step 1. Eliminate the start variable from the right-hand side of the rules.

- New start variable  $S_0$
- New rule  $S_0 \rightarrow S$

Step 2. Remove  $\epsilon$ -rules  $A \rightarrow \epsilon$ , where  $A \in V - \{S\}$ .

- Before:  $B \rightarrow xAy$  and  $A \rightarrow \epsilon | \dots$
- After:  $B \rightarrow xAy | xy$  and  $A \rightarrow \dots$

When removing  $A \rightarrow \epsilon$  rules, insert all new replacements:

- 
- Before:  $B \rightarrow AbA$  and  $A \rightarrow \epsilon | \dots$
  - After:  $B \rightarrow AbA | bA | Ab | b$  and  $A \rightarrow \dots$



# Converting CFG into CNF

## Transformation steps

In final,

All rules must be satisfied with  
above 3 requirements.

Step 3. Remove **unit rules**  $A \rightarrow B$ , where  $A \in V$ .

- Before:  $A \rightarrow B$  and  $B \rightarrow xCy$
- After:  $A \rightarrow xCy$  and  $B \rightarrow xCy$

Step 4. Eliminate all rules having more than two symbols on the right-hand side.

---

- Before:  $A \rightarrow B_1B_2B_3$
- After:  $A \rightarrow B_1A_1, A_1 \rightarrow B_2B_3$

Step 5. Eliminate all rules of the form  $A \rightarrow ab$ , where  $a$  and  $b$  are not both variables.

- Before:  $A \rightarrow ab$
- After:  $A \rightarrow B_1B_2, B_1 \rightarrow a, B_2 \rightarrow b$ .



# Converting CFG into CNF

## Example

Given a CFG  $G = (V, \Sigma, R, S)$ , where  $V = \{A, B\}$ ,  $\Sigma = \{0, 1\}$ , A is the start variable, and R consists of the rules:

$$\begin{aligned} A &\rightarrow BAB \mid B \mid \epsilon \\ B &\rightarrow 00 \mid \epsilon \end{aligned}$$

$\epsilon$ -rules:

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

Convert this G to CNF:

Step 1. Eliminate the start variable from the right-hand side of the rules.

$$S \rightarrow A$$

$$A \rightarrow BAB \mid B \mid \epsilon$$

$$B \rightarrow 00 \mid \epsilon$$



# Converting CFG into CNF

## Example

$$\begin{array}{l} S \rightarrow A \\ \checkmark A \rightarrow BAB \mid B \mid \epsilon \\ \checkmark B \rightarrow OO \mid \epsilon \end{array}$$

Step 2. Remove  $\epsilon$ -rules.

(1) Remove  $A \rightarrow \epsilon$ :  $S \rightarrow A, A \rightarrow BAB$

$$\left\{ \begin{array}{l} S \rightarrow A \mid \epsilon \\ A \rightarrow BAB \mid B \mid BB \\ B \rightarrow OO \mid \epsilon \end{array} \right.$$

(2) Remove  $B \rightarrow \epsilon$ :  $A \rightarrow BAB, A \rightarrow B, A \rightarrow BB$

$$S \rightarrow A \mid \epsilon$$

$$A \rightarrow BAB \mid B \mid BB \mid AB \mid BA \mid A$$

$$B \rightarrow OO$$



## Converting CFG into CNF

Example

$$\begin{array}{l} S \rightarrow A \\ A \rightarrow A \\ A \rightarrow B \end{array}$$

Step 3. Remove **unit-rules**.

$$\left\{ \begin{array}{l} S \rightarrow A | \epsilon \\ A \rightarrow BAB | B | BB | AB | BA | A \\ B \rightarrow \text{oo} \end{array} \right.$$

(1) Remove  $A \rightarrow A$ :

$$\left\{ \begin{array}{l} S \rightarrow A | \epsilon \\ A \rightarrow BAB | B | BB | AB | BA \\ B \rightarrow \text{oo} \end{array} \right.$$

(2) Remove  $S \rightarrow A$ :

$$\begin{array}{l} S \rightarrow B \\ A \rightarrow B \end{array} \quad \left\{ \begin{array}{l} S \rightarrow \epsilon | BAB | B | BB | AB | BA \\ A \rightarrow BAB | B | BB | AB | BA \\ B \rightarrow \text{oo} \end{array} \right.$$



# Converting CFG into CNF

## Example

Step 3. Remove **unit-rules**.

$$S \rightarrow \epsilon \mid BAB \mid B \mid BB \mid AB \mid BA$$

$$A \rightarrow BAB \mid B \mid BB \mid AB \mid BA$$

$$B \rightarrow OO$$

(3) Remove  $S \rightarrow B$ :

$$\left\{ \begin{array}{l} S \rightarrow \epsilon \mid BAB \mid BB \mid AB \mid BA \\ \checkmark A \rightarrow BAB \mid \underline{B} \mid BB \mid AB \mid 3A \\ B \rightarrow OO \end{array} \right.$$

(4) Remove  $\underline{A} \rightarrow \underline{B}$ :

$$S \rightarrow \epsilon \mid BAB \mid BB \mid AB \mid BA \mid OO$$

$$A \rightarrow BAB \mid BB \mid AB \mid BA$$

$$B \rightarrow OO$$



# Converting CFG into CNF

$\checkmark S \rightarrow \epsilon | BAB | BB | AB | BA | \text{oo}$

$A \rightarrow BAB \mid BB \mid AB \mid BA \mid \text{oc}$

B → OO

Step 4. Eliminate all rules having more than two symbols on the right-hand side.

(1) Remove  $S \rightarrow BAB$ : 想办法把  $BAB$  变为两个 symbol

$$\begin{array}{c}
 \text{BAB} \xrightarrow{\text{A}\downarrow} \text{BA}, \\
 S \rightarrow \epsilon | \text{BB} | \text{AB} | \text{BA} | \text{00} | \text{BA}, \\
 \underline{\text{A}} \rightarrow \underline{\text{BAB}} | \text{BB} | \text{AB} | \text{BA} | \text{00} \\
 \text{B} \rightarrow \text{00}
 \end{array}$$

assume  $A_1 \rightarrow AB$

(2) Remove A → BAB:

$S \rightarrow \epsilon \mid BB \mid AB \mid BA \mid \underline{\infty} \mid BA,$

replace  $\text{OO} \rightarrow A_3A_3$

$$A \rightarrow B_3 | AB | BA | \underline{OO} | BA_2$$

Replace  $\text{O}_2 \rightarrow \text{Ag}_2\text{Ag}$

$$\beta \rightarrow 00$$

$$A_1 \rightarrow AB$$

$$A_2 \rightarrow AB$$



# Converting CFG into CNF

## Example

Step 5. Eliminate all rules, whose right-hand side contains exactly two symbols, which are not both variables.

(1) Remove  $S \rightarrow 00$ :

$$S \rightarrow \epsilon \mid BB \mid AB \mid BA \mid BA_1 \mid A_3A_3$$

$$A \rightarrow BB \mid AB \mid BA \mid BA_2 \mid A_4A_4$$

$$\checkmark B \rightarrow \text{oo} \Rightarrow \text{replace to } B \rightarrow A_5A_5$$

$$A_1 \rightarrow AB$$

$$A_2 \rightarrow AB$$

$$A_3 \rightarrow O$$

$$A_4 \rightarrow O$$

$$A_5 \rightarrow O$$

(2) Remove  $A \rightarrow 00$ :

(3) Remove  $B \rightarrow 00$



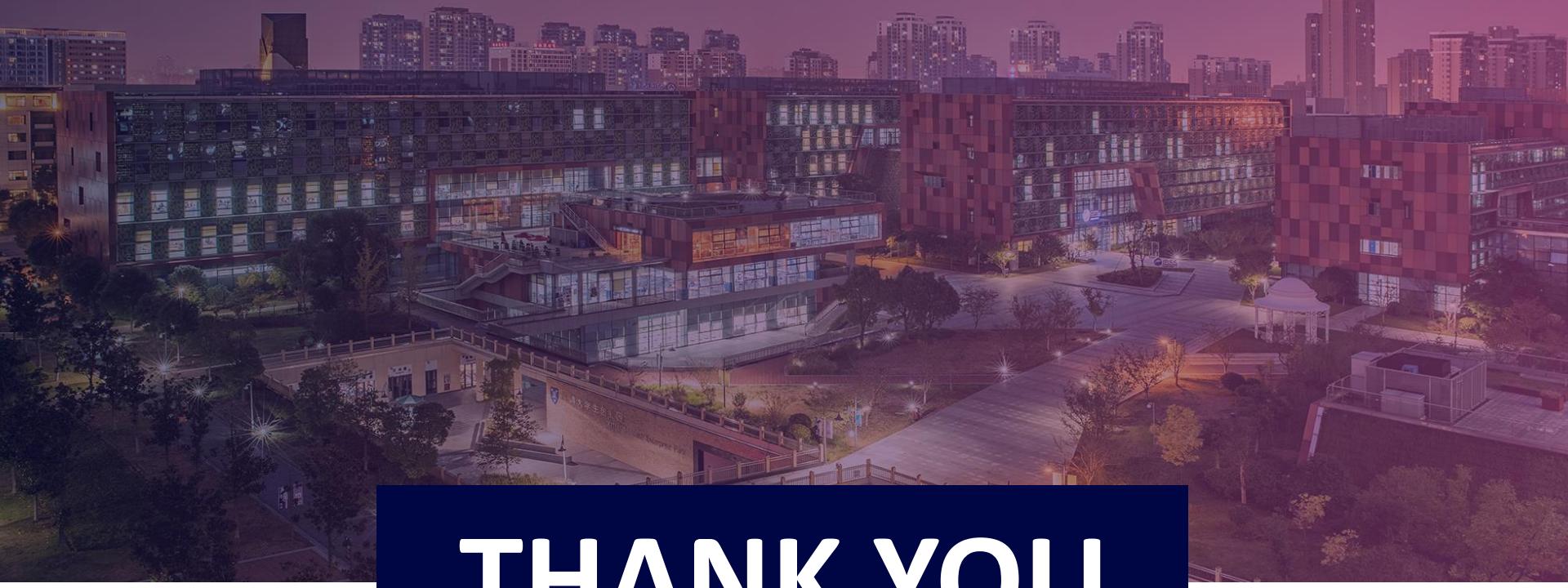
# Converting CFG into CNF

## Example

Step 5. Eliminate all rules, whose right-hand side contains exactly two symbols, which are not both variables.

(3) Remove  $S \rightarrow 00$ :





THANK YOU



Xi'an Jiaotong-Liverpool University  
西交利物浦大学

XJTLU | SCHOOL OF  
FILM AND  
TV ARTS