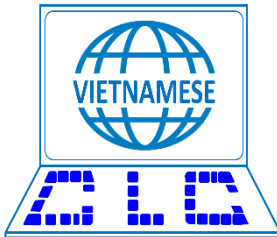


## **Section 2: Formal Languages**



**Lecturer: Assoc.Prof. Dr. Dinh Dien**

# Introduction to Formal Languages

- Formal languages are defined with respect to a given *alphabet*. The alphabet ( $\Sigma$ ) is simply a finite set of symbols. E.g.:
  - $\Sigma$  (English) = [the, a, of, in, am, I, book, teacher, mother, new, read, carry out, fall in love, ...]
  - $\Sigma$  (Vietnamese) = [của, các, cái, một, là, trong, tôi, mới, sách, giáo viên, mẹ, đọc, thực hiện, cuốn, của đảng tôi, ...]
  - $\Sigma$  (French) = [le, une, de, livre, mère, pomme de terre, ...]
  - $\Sigma$  (Chinese) = [的, 一, 书, 老师, 妈妈, 本, ...]
  - $\Sigma$  (Japanese) = [の, に, 本, お母さん, ...]
  - $\Sigma$  (Korean) = [안, 이, 그리고, 책, 선생님, 어머니, ...]

# STRING

- A finite sequence of *symbols* is called a *string*.
- The length of a string  $w$  is the number of symbols:  $|w|$
- If  $w$  is a string ( $w \in \Sigma^*$ ) and  $a$  is a symbol ( $a \in \Sigma$ ), then  $x.a$  (concatenation) will be a string on  $\Sigma$
- If  $w, y$  are strings ( $w, y \in \Sigma^*$ ), then  $x.y$  (concatenation) will be a string on  $\Sigma$  with its length:  $|xy| = |x| + |y|$ .
- Ex:  $\Sigma$  (computer language) =  $\{0,1\}$ :  $x = \text{"001"}$  is a string on  $\Sigma$  and  $a = \text{"0"}$  is a symbol of  $\Sigma$ , then  $xa = \text{"0010"}$  is a new string on  $\Sigma$ . The length:  $|x| = 3$  and  $|xa| = 4$
- But:  $b = \text{'2'} \notin \Sigma$ , then:  $xb = \text{"0012"}$  is NOT a string on  $\Sigma$ .

# SENTENCE

- $x = \text{“Tôi là một giáo_viên”}$  is a sentence on  $\Sigma$  (Vietnamese), and  $a = \text{“mới”}$  is a word of  $\Sigma$  (Vietnamese), we have a new sentence in Vietnamese:  $xa = \text{“Tôi là một giáo_viên mới”}$
- The length:  $|x| = 4$  and  $|xa| = 5$
- But: if  $b = \text{“new”} \notin \Sigma$  (Vietnamese), then  $xb = \text{“Tôi là một giáo_viên new”}$  is NOT a Vietnamese sentence.

# LANGUAGE

- Power set: from a given  $\Sigma=[0,1]$ , we have:
- $\Sigma^0=[\epsilon]$ : set of 0-symbol strings
- $\Sigma^1=[0,1]$ : set of 1-symbol strings
- $\Sigma^2=[00,01,10,11]$ : set of 2-symbol strings
- $\Sigma^3=[000,001,010,011,100,101,110,111]$
- $\Sigma^n$ : set of n-symbol strings
- $\Sigma^+$ : set of n-symbol strings with  $n \geq 1$   
 $\Sigma^+ = \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \dots \cup \Sigma^n = [0,1,00,01,10,11,000,001, \dots, 111, \dots]$
- $\Sigma^*$ : set of n-symbol strings with  $n \geq 0 \Rightarrow \Sigma^* = \Sigma^0 \cup \Sigma^+$
- The order of strings in the power set: *normalization* or dictionary (alphabet).

# LANGUAGE

- In NL, there is no  $\Sigma^0$ . Why?
- $\Sigma^1 = \{ \text{“Cháy!”}, \text{“Cướp!”}, \dots \}, \{ \text{“Fire!”}, \dots \}$  (limited)
- $\Sigma^2 = \{ \text{“Tôi đi”}, \text{“Anh ngủ”}, \text{“Trời mưa”}, \text{“Giáo_viên nghiên_cứu”}, \text{“Sinh_viên thực_hiện”}, \dots \}, \{ \text{“It rains”}, \text{“I go”}, \text{“You sleep”}, \text{“Teachers research”}, \text{“Students carry_out”}, \dots \}$  (many)
- $\Sigma^3, \Sigma^4, \dots, \Sigma^n$  : increase exponentially
- Ex:  $\Sigma(\text{computer}) = \{0, 1\}$
- $\Rightarrow \Sigma^3 = [000, 001, 010, 011, 100, 101, 110, 111]$
- ✓ Only *italic* strings are valid.
- $L = [001, 011, 100] \subset \Sigma^3 \subset \Sigma^*$

# LANGUAGE DEFINITIONS

- DEFINITION OF A LANGUAGE: is a *subset* of  $\Sigma^*$
- How to determine that *subset*?
- If the number is *limited*, we can *enumerate* them.
- How large is the subset ? Limited or Unlimited ?
- For  $\Sigma=[0,1]$  (its vocabulary has 2 words only)
- $\Sigma^8$ :  $2^8 = 256$  cases;  $\Sigma^{16}$  :  $2^{16} = 65.536$  ;  $\Sigma^{32}$  :  $2^{32} \sim 4,2$  billion.  $\Sigma^{64}$ :  $2^{64} \sim 18$  million billion (quintillion) cases.
- The size of this “rice heap”(chess story)is 32 billion  $m^3$ , approx.12,000 times of the biggest pyramid (Khufu,Egypt)
- Luckily, the number of valid strings (instruction) is small/limited (opcode)  $\Rightarrow$  instruction set of CPU.
- But: in other languages ?

# LANGUAGE PRINCIPLES

- $\Sigma=[0,1]$ ,  $L_1= [11,011,101,0101,0110,...]$ ,  $L_1 \subset \Sigma^*$
- $L_1$ 's principle? *The number of '1s' is 2.*
- $L_2= [11,011,0101,0110,11011,110011,1011111,1101111,...]$
- $L_2$ 's principle? *The number of '1s' is even.*
- $L_3= [01,0101,0110, 010011, 011100, 01110010,10010101,...]$
- $L_3$ 's principle? *The number of '1s' = '0s'.*
- $L_1, L_2, L_3$ : infinite countable sets
  - *What is the principle of NL?*
  - *The principles of phonetics/phonology, morphology, grammar(syntax), SEMANTICS, pragmatics, ...*



# NATURAL LANGUAGES

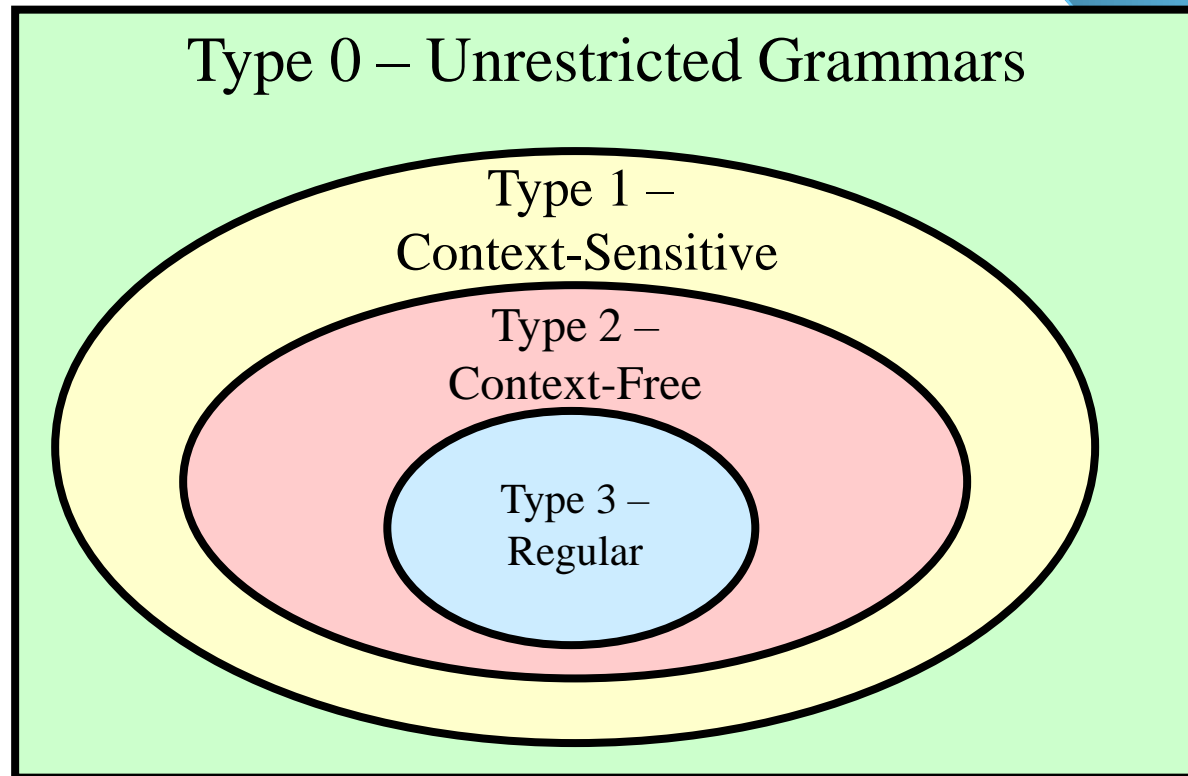
- $\Sigma^+ =$  [“Tôi là một”, “Tôi là một mới”, “là giáo\_viên một mới sách”, “Tôi đọc giáo\_viên sách mới”, “*Tôi đọc sách giáo\_viên mới*”, “*Mẹ đọc sách giáo\_viên mới*”, “Đọc mới giáo\_viên sách là”, “*Mẹ mới đọc sách giáo\_viên*”, “của sách một”,...]
- Italic sentences: grammatically & semantically correct.
- The number of correct sentences: huge!
- $\Sigma$  (Vietnamese)  $\sim$  40k words, the avg length  $\sim$  20w
- Candidate sentences:  $40.000^{20}$ ;  $40.000 > 2^{15}$  ( $= 32.768$ )  $\Rightarrow 40.000^{20} > (2^{15})^{20} = 2^{300} = 2^{30 \times 10} = (10^9)^{10} = (10^{90})$  (90 ‘0s’ !)
- The number of sand grains on the Earth:  $10^{21}$
- The number of stars in the universe:  $7 \times 10^{22}$  (sextillion)

# NATURAL LANGUAGES

- Possible Vietnamese sentences:  $10^{40}$
- English (Robert Mannell at Macquarie, Aus):  $10^{570}$   
[http://clas.mq.edu.au/speech/infinite\\_sentences/index.html](http://clas.mq.edu.au/speech/infinite_sentences/index.html)
- Extremely huge!
- We need another way to determine the subset of NL
- That is: GRAMMAR
- SUBSET = {collection of grammatically correct sentences}

# Types of Grammars: Chomsky hierarchy

- 0. UG:  $\alpha \rightarrow \beta$  with  $\beta \in (N \cup \Sigma)^*$  and  $\alpha \in (N \cup \Sigma)^+$
- 1. CSG:  $\gamma_1 X \gamma_2 \rightarrow \gamma_1 \alpha \gamma_2$  with  $X, \alpha, \gamma_1, \gamma_2 \in (N \cup \Sigma)^+$
- 2. CFG:  $X \rightarrow \alpha$  with  $X \in N, \alpha \in (N \cup \Sigma)^+$
- 3. RG:  $X \rightarrow a$  and  $X \rightarrow aY$  or  $X \rightarrow Yb$  with  $X, Y \in N, a, b \in \Sigma$



# DERIVATION

- Let  $G(V, \Sigma, P, S)$  be a CFG.
- Let  $w_0 = \alpha X \beta$  (the concatenation of  $\alpha, X$  and  $\beta$ ) and  $w_1 = \alpha \gamma \beta$  be strings over  $\Sigma$ .
- If  $X \rightarrow \gamma$  is a production of  $P$  we say that  $w_1$  is **directly derivable** from  $w_0$  and we write  $w_0 \Rightarrow w_1$
- If  $w_0, w_1, \dots, w_n$  are strings over  $V$  such that  $w_0 \Rightarrow w_1, w_1 \Rightarrow w_2, \dots, w_{n-1} \Rightarrow w_n$ , then we say that  $w_n$  is derivable from  $w_0$ , and write  $w_0 \Rightarrow^* w_n$
- The sequence of steps used to obtain  $w_n$  from  $w_0$  is called a **derivation**.

# DEFINITION OF LANGUAGES

- Let  $G(V, \Sigma, P, S)$  be a CFG.
- The language generated by  $G$  denoted by  $L(G)$ , is the set of all strings of terminals that are derivable from the start symbol  $S$ .
- $$L(G) = \{w \mid (w \in \Sigma^*) \cap (S \Rightarrow^* w)\}$$

# Derivation Order

Consider the following example grammar with 5 productions:

- |                       |                            |                            |
|-----------------------|----------------------------|----------------------------|
| 1. $S \rightarrow AB$ | 2. $A \rightarrow aaA$     | 4. $B \rightarrow Bb$      |
|                       | 3. $A \rightarrow \lambda$ | 5. $B \rightarrow \lambda$ |

- |                       |                            |                            |
|-----------------------|----------------------------|----------------------------|
| 1. $S \rightarrow AB$ | 2. $A \rightarrow aaA$     | 4. $B \rightarrow Bb$      |
|                       | 3. $A \rightarrow \lambda$ | 5. $B \rightarrow \lambda$ |

Leftmost derivation order of string :

$$\begin{array}{ccccccccc} 1 & & 2 & & 3 & & 4 & & 5 \\ S & \Rightarrow & AB & \Rightarrow & aaAB & \Rightarrow & aaB & \Rightarrow & aaBb & \Rightarrow & aab \end{array}$$

At each step, we substitute the leftmost variable

- |                       |                            |                            |
|-----------------------|----------------------------|----------------------------|
| 1. $S \rightarrow AB$ | 2. $A \rightarrow aaA$     | 4. $B \rightarrow Bb$      |
|                       | 3. $A \rightarrow \lambda$ | 5. $B \rightarrow \lambda$ |

Rightmost derivation order of string :  $aab$

$$\begin{array}{ccccccccc} & 1 & & 4 & & 5 & & 2 & & 3 \\ S & \Rightarrow & AB & \Rightarrow & ABb & \Rightarrow & Ab & \Rightarrow & aaAb & \Rightarrow & aab \end{array}$$

At each step, we substitute the  
rightmost variable



- |                       |                            |                            |
|-----------------------|----------------------------|----------------------------|
| 1. $S \rightarrow AB$ | 2. $A \rightarrow aaA$     | 4. $B \rightarrow Bb$      |
|                       | 3. $A \rightarrow \lambda$ | 5. $B \rightarrow \lambda$ |

Leftmost derivation of :  $aab$

$$\begin{array}{ccccccccc}
 1 & & 2 & & 3 & & 4 & & 5 \\
 S & \Rightarrow & AB & \Rightarrow & aaAB & \Rightarrow & aaB & \Rightarrow & aaBb & \Rightarrow & aab
 \end{array}$$

Rightmost derivation of :  $aab$

$$\begin{array}{ccccccccc}
 1 & & 4 & & 5 & & 2 & & 3 \\
 S & \Rightarrow & AB & \Rightarrow & ABb & \Rightarrow & Ab & \Rightarrow & aaAb & \Rightarrow & aab
 \end{array}$$

# Derivation Trees

Consider the same example grammar:

$$S \rightarrow AB \quad A \rightarrow aaA \mid \lambda \quad B \rightarrow Bb \mid \lambda$$

And a derivation of  $aaab$ :

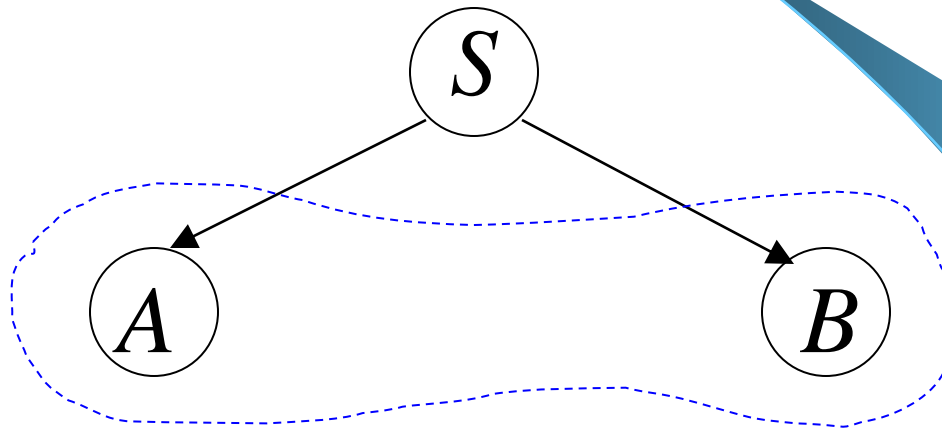
$$S \Rightarrow AB \Rightarrow aaAB \Rightarrow aaABb \Rightarrow aaBb \Rightarrow aaab$$

$$S \rightarrow AB$$

$$A \rightarrow aaA \mid \lambda$$

$$B \rightarrow Bb \mid \lambda$$

$$S \Rightarrow AB$$



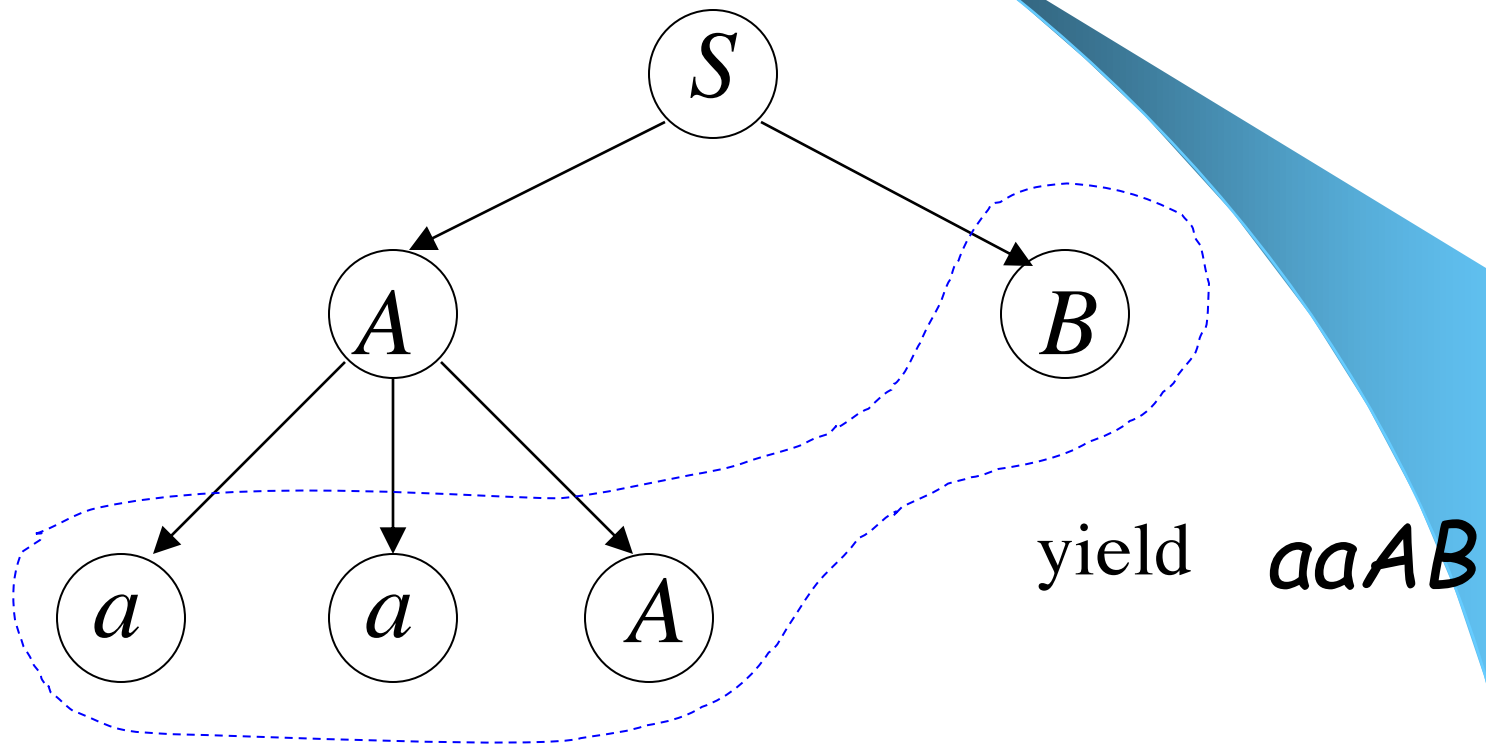
yield ***AB***

$$S \rightarrow AB$$

$$A \rightarrow aaA \mid \lambda$$

$$B \rightarrow Bb \mid \lambda$$

$$S \Rightarrow AB \Rightarrow aaAB$$

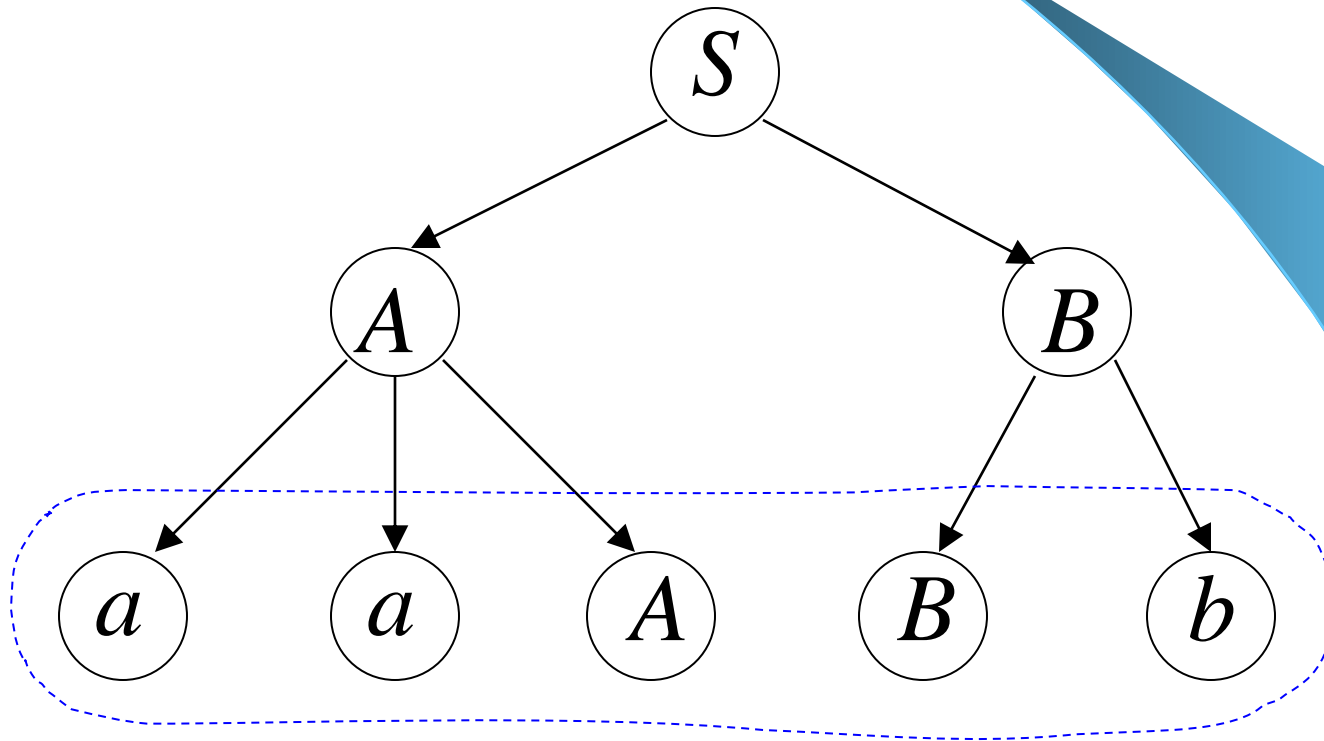


$$S \rightarrow AB$$

$$A \rightarrow aaA \mid \lambda$$

$$B \rightarrow Bb \mid \lambda$$

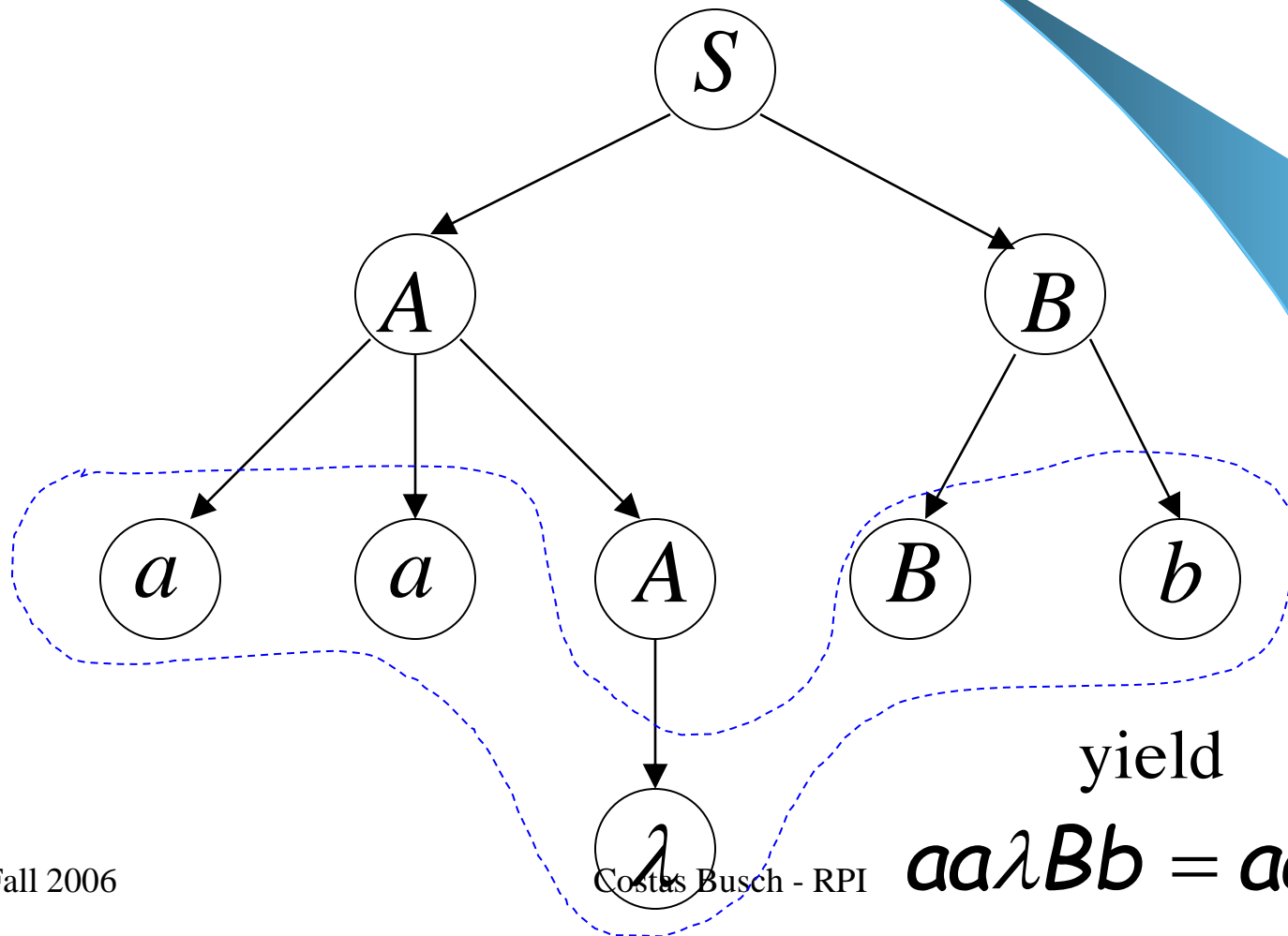
$$S \Rightarrow AB \Rightarrow aaAB \Rightarrow aaABb$$



yield  $aaABb$

$$S \rightarrow AB \quad A \rightarrow aaA \mid \lambda \quad B \rightarrow Bb \mid \lambda$$

$$S \Rightarrow AB \Rightarrow aaAB \Rightarrow aaABb \Rightarrow aaBb$$



$$aa\lambda Bb = aaBb$$

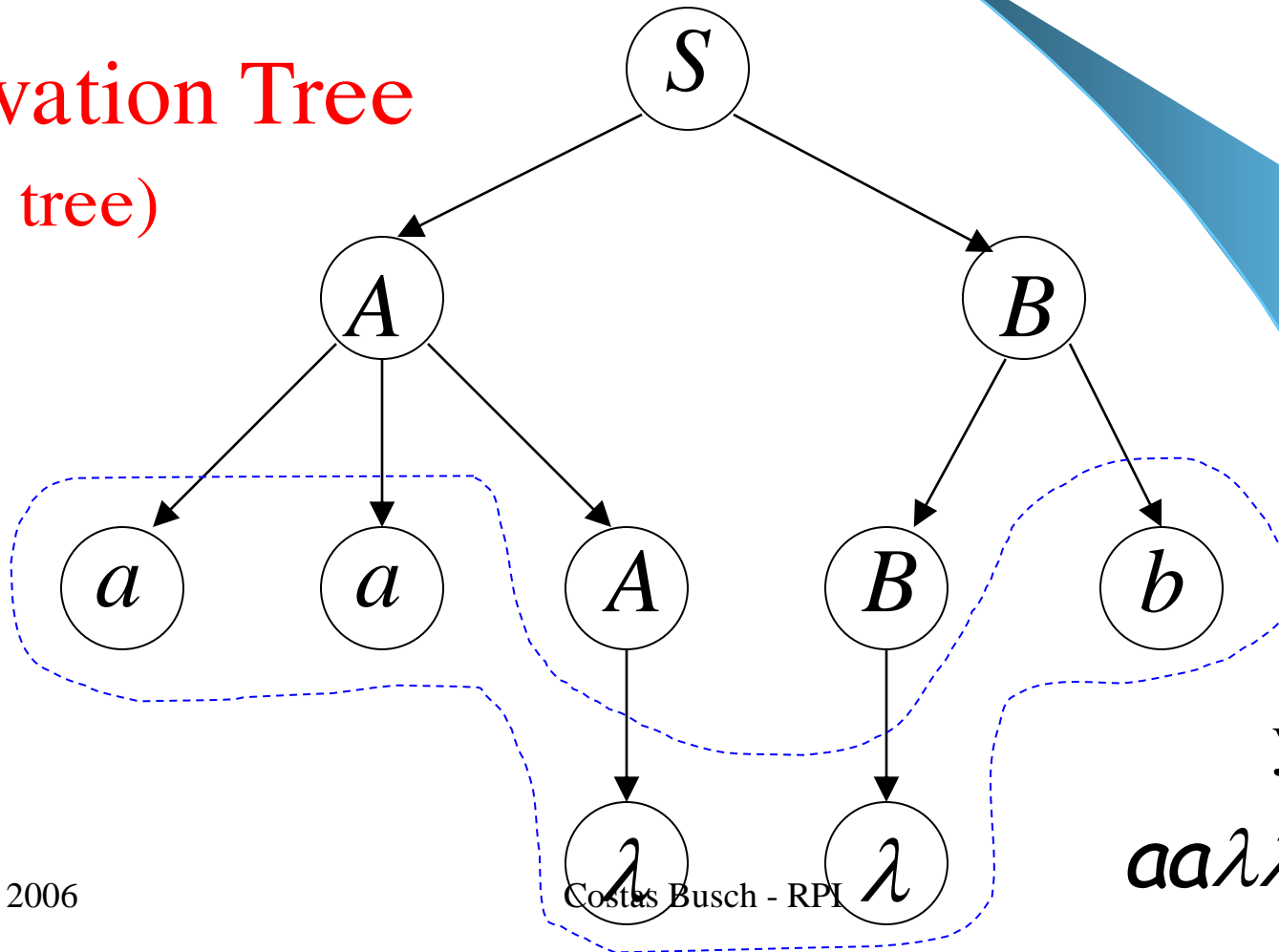
$$S \rightarrow AB$$

$$A \rightarrow aaA \mid \lambda$$

$$B \rightarrow Bb \mid \lambda$$

$$S \Rightarrow AB \Rightarrow aaAB \Rightarrow aaABb \Rightarrow aaBb \Rightarrow aab$$

Derivation Tree  
(parse tree)

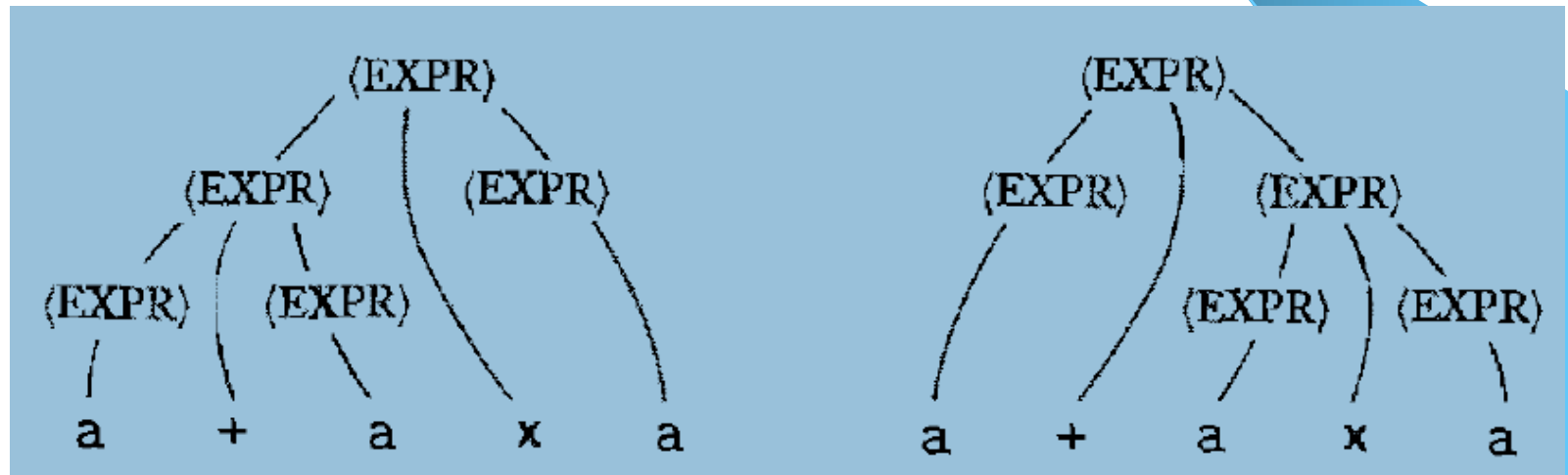


yield

$$aa\lambda\lambda b = aab$$

# AMBIGUOUS

- ❑ A sentence which can derive more than one parse *tree* is **ambiguous**.
- ❑ Example: a grammar for arithmetic expressions:  
 $\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle \mid \langle \text{EXPR} \rangle . \langle \text{EXPR} \rangle \mid ( \langle \text{EXPR} \rangle ) \mid a$   
✓ String:  $a + a.a$



- ❑ A language in which contains at least one ambiguous sentence is **ambiguous**.
- Give examples of ambiguous sentences in NL?



+ *Equivalent* grammar:

G1 and G2 are equivalent iff  $L(G1) = L(G2)$

+ *Recursive*:

Given a  $G = (N, \Sigma, P, S)$ , a non-term  $X$  of  $G$  is recursive if there are following production rules:

$X \Rightarrow \alpha X\beta$  with  $\alpha, \beta \in (\Sigma \cup N)^*$

If  $\alpha = \varepsilon$  : left recursive

If  $\beta = \varepsilon$  : right recursive

If  $\alpha, \beta \neq \varepsilon$ : middle recursive

Indirect recursive ?

Give examples of recursive structures in NL?

+Categorized grammar

Given a grammar  $G = (N, \Sigma, T, P, S)$ .

$T$  is set of categories of terminals.

$$\Sigma = \cup T_i$$

With  $i \neq j$ :  $T_i \cap T_j = \emptyset$ ?

Give examples in English, Vietnamese, etc.

+ *CNF* (Chomsky Normal Form):if every production rules in  $P$  have following forms:

$$X \rightarrow YZ \text{ with } X, Y, Z \in N, \text{ or}$$

$$X \rightarrow a \text{ with } a \in \Sigma$$

Can CFG be converted into CNF ?

# Backus-Naur Form

- $\langle \text{sentence} \rangle ::= \langle \text{noun phrase} \rangle \langle \text{verb phrase} \rangle$
- $\langle \text{noun phrase} \rangle ::= \langle \text{article} \rangle [\langle \text{adjective} \rangle] \langle \text{noun} \rangle$
- $\langle \text{verb phrase} \rangle ::= \langle \text{verb} \rangle [\langle \text{adverb} \rangle]$
- $\langle \text{article} \rangle ::= \text{a} / \text{the}$
- $\langle \text{adjective} \rangle ::= \text{large} | \text{hungry}$
- $\langle \text{noun} \rangle ::= \text{rabbit} | \text{mathematician}$
- $\langle \text{verb} \rangle ::= \text{eats} / \text{hops}$
- $\langle \text{adverb} \rangle ::= \text{quickly} / \text{wildly}$

Square brackets []  
mean “optional”

Vertical bars  
mean “alternatives”

Given a grammar:

$S ::= NP VP$

$NP ::= det noun$

$NP ::= det adj noun$

$NP ::= det noun PP$

$VP ::= verb$

$PP ::= prep NP$

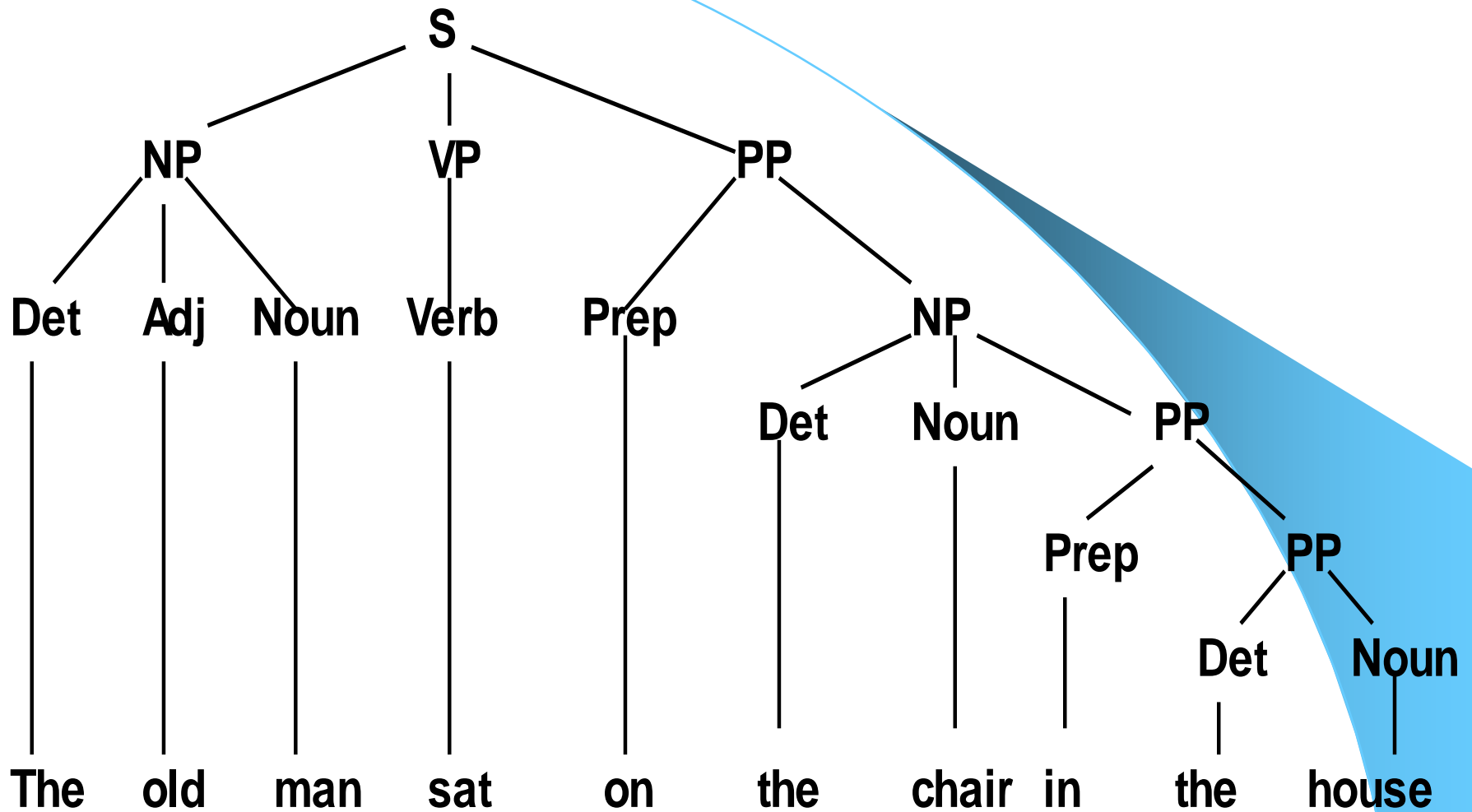
$PP ::= prep NP PP$

Is that CFG ?

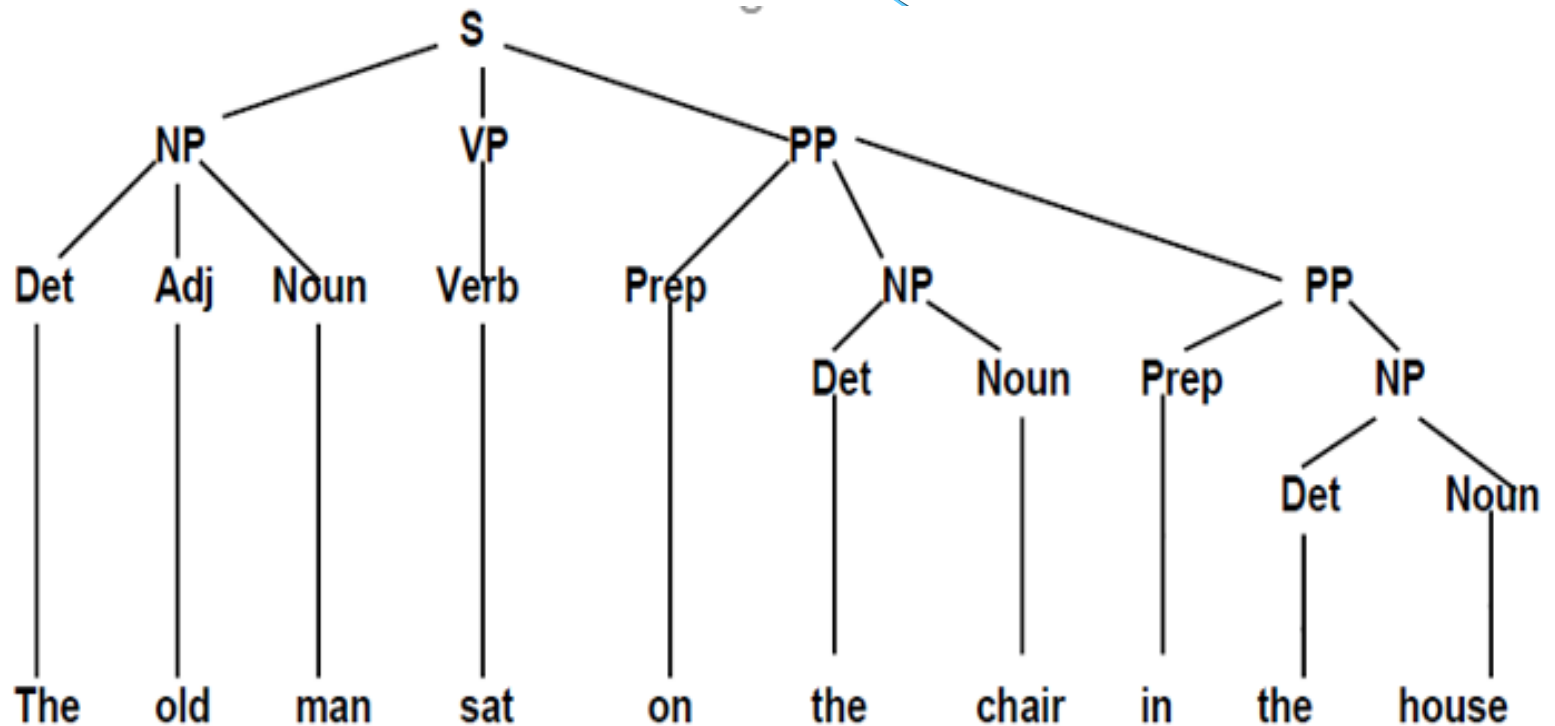
Parse: “The old man sat on the chair in the house”

Is that sentence ambiguous?

# Derivation tree #1



# Derivation tree #2



# EARLEY PARSING ALGORITHM

Given a CFG grammar:

$S \rightarrow NP VP$

$S \rightarrow NP VP PP$

$NP \rightarrow d NP3$

$NP3 \rightarrow a NP3$

$NP3 \rightarrow n$

$NP3 \rightarrow n PP$

$PP \rightarrow p NP2$

$NP2 \rightarrow d NP3$

$VP \rightarrow v$

- In which: d (det), a (adj), p (prep), n (noun), v (verb): are terminals
- S, NP, NP3, NP2, VP, PREPS : non-terminal symbols
- Parse the sentence: “ The young student sat in the class”
- POS: “The/d young/a student/n sat/v in/p the/d class/n”
- Formally:  $w = a_1 a_2 a_3 a_4 a_5 a_6 a_7$
- With:  $a_1=d$ ;  $a_2=a$ ;  $a_3=n$ ;  $a_4=v$ ;  $a_5=p$ ;  $a_6=d$ ;  $a_7=n$ ;



# Building the parsing list

- (1) If  $S \rightarrow \alpha \in P$ , then add  $[S \rightarrow \bullet \alpha, 0]$  into  $I_0$
- Repeat (2) and (3) until no new added-item in  $I_0$
- (2) If  $[B \rightarrow \gamma \bullet, 0] \in I_0$ , then add:  
 $[A \rightarrow \alpha B \bullet \beta, 0]$  for all items  $[A \rightarrow \alpha \bullet B \beta, 0]$  in  $I_0$ .
- (3) If  $[A \rightarrow \alpha \bullet B \beta, 0] \in I_0$ , add into  $I_0$ , all items  $[B \rightarrow \bullet \gamma, 0]$  (if no exist in  $I_0$ ).

# Building the parsing list

- Building the next table:  $I_j$
- (4) For each item  $[B \rightarrow \alpha \bullet a \beta, i]$  in table  $I_{j-1}$  with  $a = a_j$ , we add item  $[B \rightarrow \alpha a \bullet \beta, i]$  into table  $I_j$
- Repeat (5) and (6) until no new added-item in  $I_j$ .
- (5) If  $[A \rightarrow \alpha \bullet B\beta, i] \in I_j$ , we add into  $I_j$  items  $[B \rightarrow \bullet \gamma, j]$ .
- (6) If  $[A \rightarrow \alpha \bullet, i] \in I_j$ , search in  $I_i$  for items  $[B \rightarrow \alpha \bullet A\beta, k]$
- If found, we add item  $[B \rightarrow \alpha A \bullet \beta, k]$  into table  $I_j$ .

- Table I0:
- $S \rightarrow \bullet NP VP, 0$
- $S \rightarrow \bullet NP VP PP, 0$
- $NP \rightarrow \bullet d NP3, 0$
- Table I1: (The/ $a_1=d$ )
- $NP \rightarrow d \bullet NP3, 0$
- $NP3 \rightarrow \bullet a NP3, 1$
- $NP3 \rightarrow \bullet n, 1$
- $NP3 \rightarrow \bullet n PP, 1$
- Table I2 : (young/  $a_2=a$ )
- $NP3 \rightarrow a \bullet NP3, 1$
- $NP3 \rightarrow \bullet a NP3, 2$
- $NP3 \rightarrow \bullet n, 2$
- $NP3 \rightarrow \bullet n PP, 2$

○ Table I3 : (student/  $a_3=n$ )

- $NP3 \rightarrow N \bullet, 2$
- $NP3 \rightarrow N \bullet PP, 2$
- $PP \rightarrow \bullet p NP2, 3$
- $NP3 \rightarrow a NP3 \bullet, 1$
- $NP \rightarrow d NP3 \bullet, 0$
- $S \rightarrow NP \bullet VP, 0$
- $S \rightarrow NP \bullet VP PP, 0$
- $VP \rightarrow \bullet v, 3$

○ Table I4 : (sat/  $a_4=v$ )

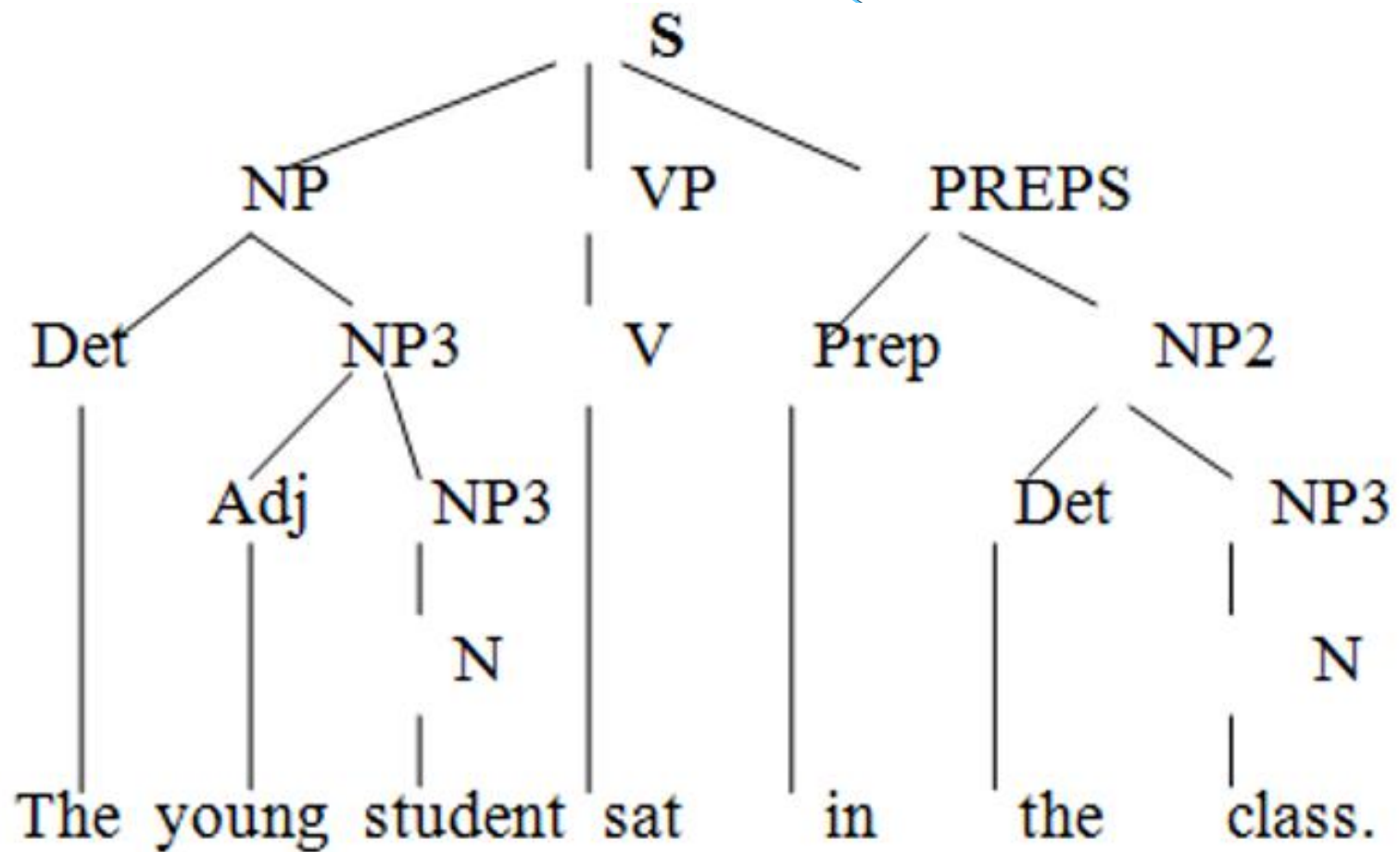
- $VP \rightarrow V \bullet, 3$
- $S \rightarrow NP VP \bullet, 0 *$
- $S \rightarrow NP VP \bullet PP, 0$
- $PP \rightarrow \bullet p NP2, 4$

- Table I5 : (in/a<sub>5</sub>=p)
- PP → p • NP2, 4
- NP2 → • d NP3, 5
- Table I6 : (the/a<sub>6</sub>=d)
- NP2 → d • NP3, 5
- NP3 → • a NP3, 6
- NP3 → • n, 6
- NP3 → • n PP, 6
- Table I7 : (class/a<sub>7</sub>=n)
- NP3 → n •, 6
- NP3 → n • PP, 6
- PP → • p NP2, 7
- NP2 → d NP3 •, 5
- PP → p NP2 •, 4
- S → NP VP PP •, 0 \*\*

# ALGORITHM FOR BUILDING THE PARSING TREE

- From list of tables  $\Rightarrow$  building the parsing tree
- If  $[S \rightarrow \alpha \bullet, 0] \notin I_n$  (last table)  $\Rightarrow$  wrong sentence  $\Rightarrow$  no tree
- Else: correct sentence  $\Rightarrow$  building the tree:
- Initialize a global variable:  $k = n$  (length of  $w$ )
- (1) Assign:  $A = S$  and  $\alpha = X_1 X_2 \dots X_m$ . From  $R \rightarrow L$ :
- (2) If  $X_i \in \Sigma$ ,  $k--$
- (3) If  $X_i \in N$ , find item  $[X_i \rightarrow \gamma \bullet, r]$  in table  $I_k$
- (4) Assign  $A = X_i$  and  $\gamma = X_1 X_2 \dots X_m$ .
- (5) Repeat from (2) until  $k=0$

# SYNTAX TREE



- **Homework:** using above-mentioned grammar to:
- parse the sentence: “An old man sat on the new chair in the house”
- If Ok  $\Rightarrow$  draw the parsing tree.
- Else: modify the grammar to accept that sentence.
- Note: avoid “flip-flop” phenomenon.