

# CS & IT ENGINEERING

## Compiler Design

Lexical Analysis & Syntax Analysis



### Lecture No. 4



By- DEVA Sir

## TOPICS TO BE COVERED

H.W. Questions

Basics of Parsing

→ CFG

→ Definitions

→ Types of Derivations

→ Ambiguous & unambiguous

⑯

 $\text{for}(1, 2, 3); \Rightarrow \text{Syntax Error}$ 

⑰

 $\text{f80}(1, 2, 3); \Rightarrow \text{No Syntax Error}$   
[either semantic error  
or linker error]  
H.W.

⑱

 $\text{for}(2; 2; 2); \Rightarrow \text{No compilation error}$ 

⑲

 $\text{f80}(2; 2; 2); \Rightarrow \text{Syntax error}$

```
#include <stdio.h>
```

```
Void main( )
```

```
{
```

```
f();
```

```
}
```

Linked Form

Compiler  
assumes  
 $f()$  definition  
may exist in other file

```
Void main( )
```

```
{
```

```
f();
```

```
}
```

Semantic Form

# MCQ

Consider the following ANSI C program:

```
int main() {  
    Integer x;  
    Return 0;  
}
```

→ Semantic Err  
→ No syntax Err

typedef int Integer;

is missing

Which one of the following phases in a seven-phase C compiler will throw an error?

[GATE-2021-CS: 1M]

- Lexical analyzer
- Syntax analyzer
- Semantic analyzer
- Machine dependent optimizer

# MCQ

Consider the following statements:

- I. Symbol table is accessed only during lexical analysis and syntax analysis. FALSE
- II. Compilers for programming languages that support recursion necessarily need heap storage for memory allocation in the run-time environment. FALSE
- III. Errors violating the condition (any variable must be declared before its use) are detected during syntax analysis. FALSE

Which of the above statements is/are TRUE?

[GATE-2020-CS: 1M]

- I only
- I and III only
- II only
- None of I, II and III

**MCQ**

A lexical analyzer uses the following patterns to recognize three tokens  $T_1, T_2$  and  $T_3$  over the alphabet  $\{a, b, c\}$ .

$$x? = \epsilon + x$$

$$\boxed{\epsilon + x}$$

Note that ' $x?$ ' means 0 or 1 occurrence of the symbol  $x$ . Note also that the analyzer outputs the token that matches the longest possible prefix.

If the string  $bbaacabc$  is processed by the analyzer, which one of the following is the sequence of tokens it outputs?

$T_1 T_2 T_3$

$T_1 T_1 T_3$

$T_2 T_1 T_3$

$T_3 T_3 T_3 T_3$

$$\begin{aligned}
 T_1 &: a? (b|c)^* a = (\epsilon + a) (\underline{b} + \underline{c})^* \underline{a} \\
 T_2 &: b? (a|c)^* b = (\epsilon + b) (\underline{a} + \underline{c})^* \underline{b} \\
 T_3 &: c? (b|a)^* c = (\epsilon + c) (\underline{b} + \underline{a})^* \underline{c}
 \end{aligned}$$

|       | bbaacabc |   |   |   |
|-------|----------|---|---|---|
|       |          |   |   |   |
|       |          |   |   |   |
| $T_1$ | ✓        | ✓ | ✓ | X |
| $T_2$ | ✓        | ✓ | X |   |
| $T_3$ | ✓        | ✓ | ✓ | ✓ |

[GATE-2018-CS: 2M]

**MCQ**

Match the following according to input from List-I to the compiler phase in the List-II that processes it:

List-I

(P) Syntax tree

(Q) Character stream

(R) Intermediate representation

(S) Token stream

List-II      **phase**

(i) Code generator

(ii) Syntax analyser

(iii) Semantic analyser

(iv) Lexical analyser [GATE-2017-CS: 1M]

A

 $P \rightarrow (ii), Q \rightarrow (iii), R \rightarrow (iv), S \rightarrow (i)$ 

B

 $P \rightarrow (ii), Q \rightarrow (i), R \rightarrow (iii), S \rightarrow (iv)$ 

C

 $P \rightarrow (iii), Q \rightarrow (iv), R \rightarrow (i), S \rightarrow (ii)$ 

D

 $P \rightarrow (i), Q \rightarrow (iv), R \rightarrow (ii), S \rightarrow (iii)$ 

I.C.

Parse tree

(after changing  
given data)

|           |
|-----------|
| $P - i$   |
| $Q - iv$  |
| $R - i$   |
| $S - iii$ |

Given data

# MCQ

Match the following

List-I

P: Lexical analysis

Q: Top down parsing

R: Semantic analysis

S: Runtime environment

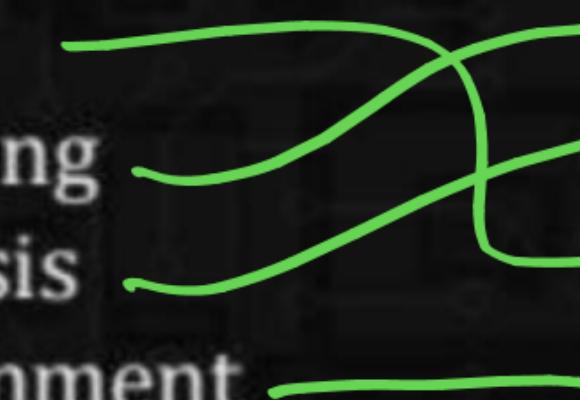
List-II

(i) Leftmost derivation

(ii) Type checking

(iii) Regular expressions

(iv) Activation records



[GATE-2016-CS: 1M]

P→(i), Q→(ii), R→(iv), S→(iii)

P→(iii), Q→(i), R→(ii), S→(iv)

P→(ii), Q→(iii), R→(i), S→(iv)

P→(iv), Q→(i), R→(ii), S→(iii)

**MCQ****TOKENS**

In a compiler, **keywords** of a language are recognized during

[GATE-2021-CS: 1M]

- parsing of the program (syntax)
- the code generation (Assembly code generate)
- the lexical analysis of the program
- dataflow analysis → Code optimization

# MCQ

Which data structure in a compiler is used for managing **information** about **variables** and **their attributes**?  
**Identifiers**

[GATE-2021-CS: 1M]

- Abstract syntax tree
- Symbol table
- Semantic stack
- Parse table

Context Free  
(Structure)  
(Syntax)

S + V + O

I am student



I are student

Context Sensitive  
(Semantic)  
(Meaning)

I am student (semantically correct)

I are student (semantically wrong)

Context Free  
(Structure)  
(Syntax)

Integer x ;

Syntax correct

Context Sensitive  
(Semantic)  
(Meaning)

Integer x ;

Semantically wrong

# Context Free Grammars [CFG]



→ It can be used to represent syntax/structure of programming languages.

$$CFG = (V, T, P, S)$$

Each rule in  $P$ :

LHS  $\rightarrow$  RHS

$$\boxed{V \rightarrow \text{Any sequence}}$$
$$V \rightarrow (VUT)^*$$

→ Set of Variables (non-terminals)  
→ Set of terminals  
→ Set of Rules (productions) (production Rule)  
→ start symbol  $S \in V$

Identify CFGs:

①  $S \rightarrow \epsilon$



②  $S \rightarrow Sab | a$



③  $S \rightarrow ABS | a$

$\textcircled{AB} \rightarrow a$



④  $S \rightarrow aSb | SS | a$



⑤  $S \rightarrow S | a | \epsilon$



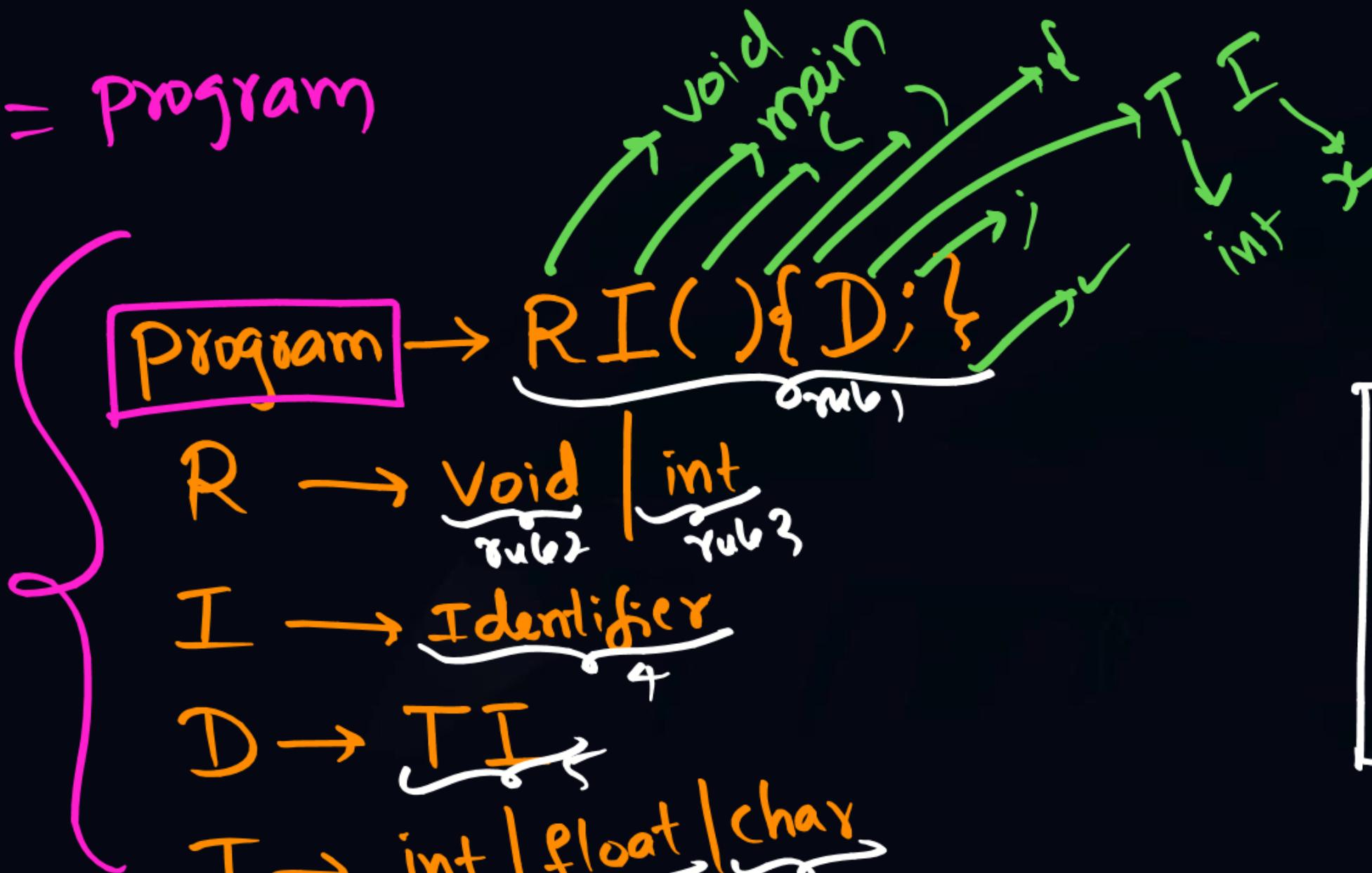
⑥  $S \rightarrow aSb$

$\textcircled{aS} \rightarrow bb$



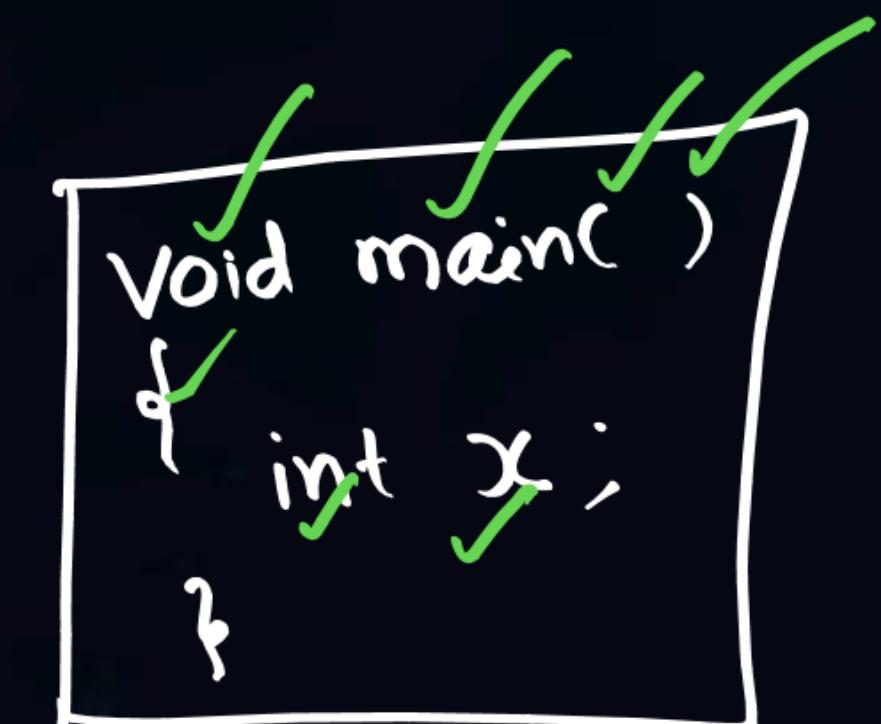
$S = \text{program}$

8 rules



$V = \{ \text{program}, R, I, D, T \}$

$T = \{ (, ), \{, \}, ;, \}, \text{void}, \text{int}, \text{Identifier}, \text{float}, \text{char} \}$



$$S \rightarrow ab$$

OR

OR

$$S \rightarrow Ab$$
$$A \rightarrow aB$$
$$B \rightarrow \epsilon$$
$$S \rightarrow Ab$$
$$A \rightarrow a$$

OR

$$S \rightarrow A\bar{B}$$
$$A \rightarrow a$$
$$\bar{B} \rightarrow b$$
$$ab$$
  
$$\cancel{\cancel{}}$$

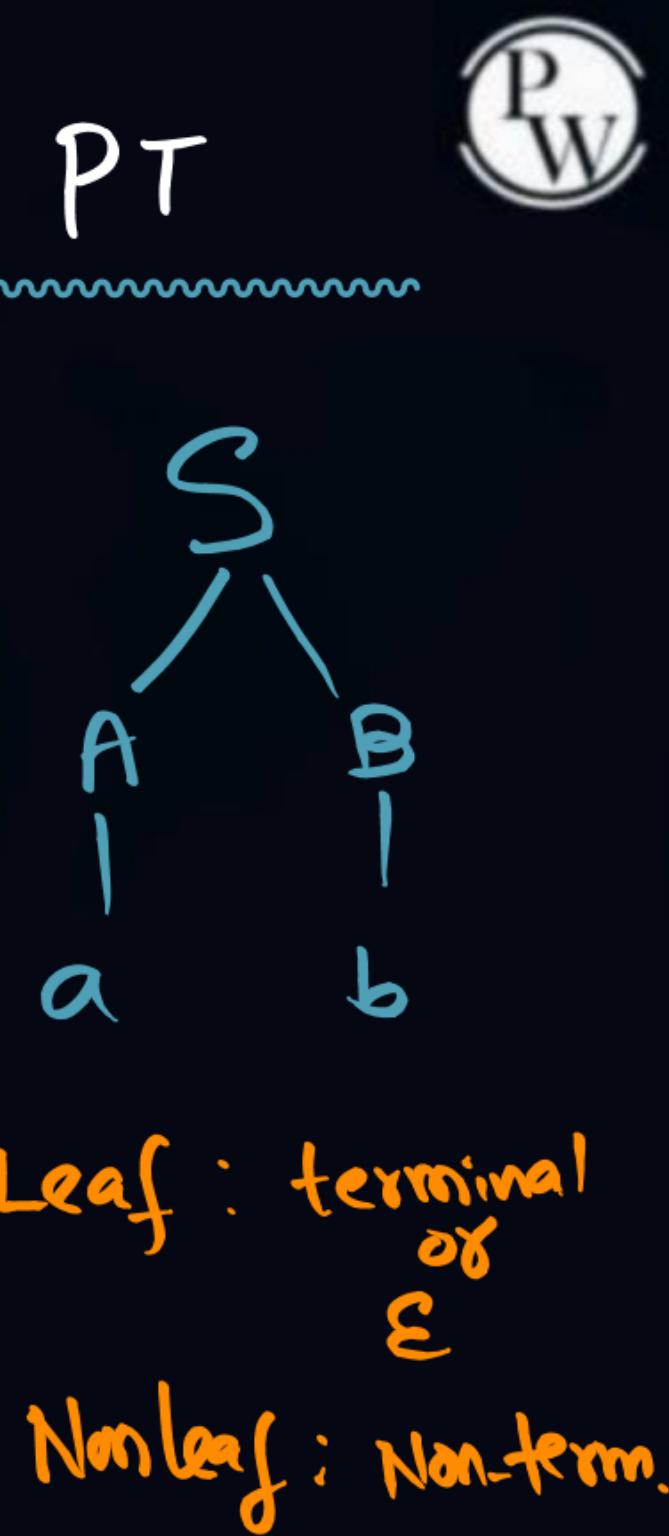
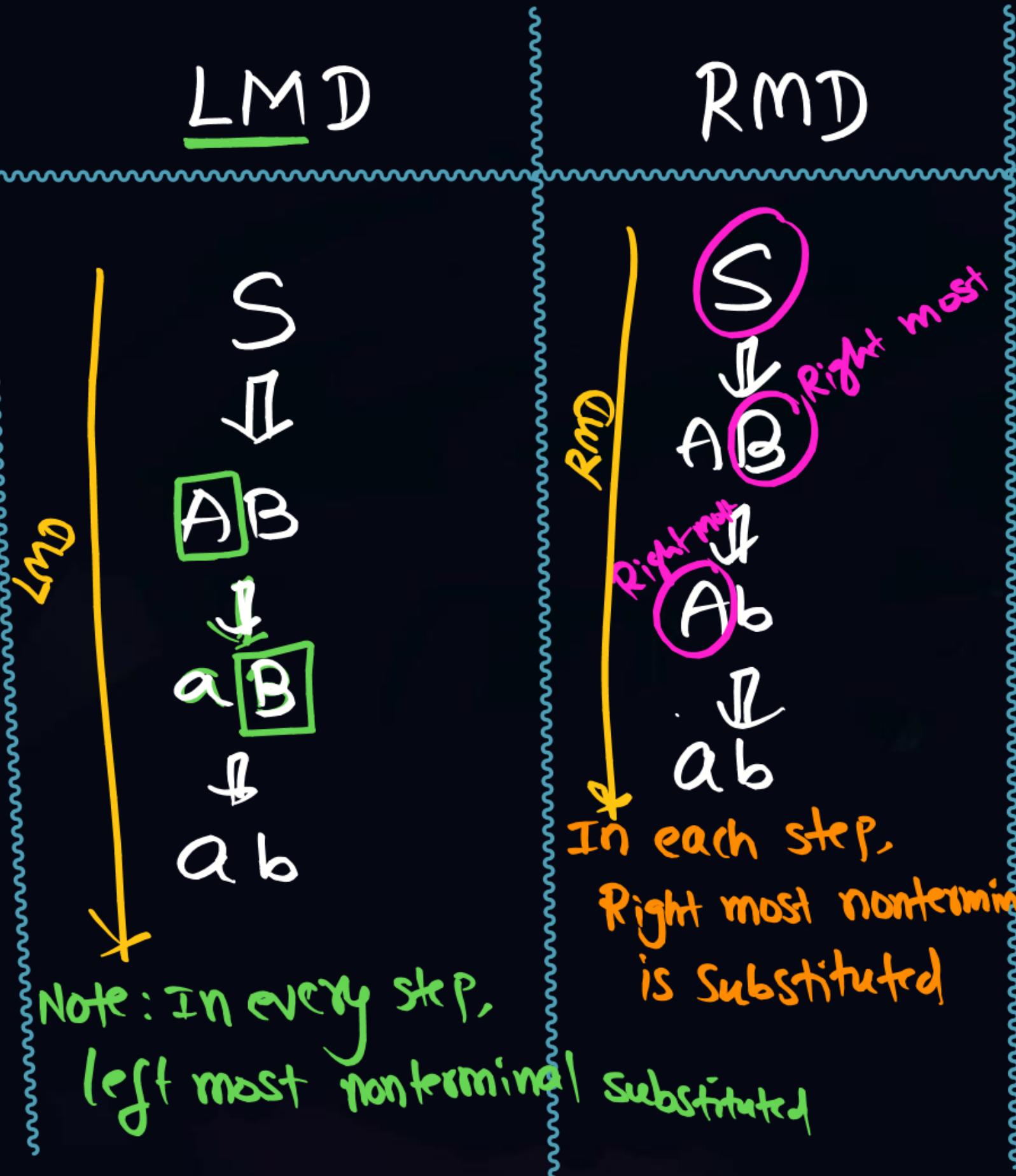
# Derivation of a String:

- Linear { I) Left MOST Derivation(LMD)  
II) Right MOST Derivation(RMD)
- Non-linear { III) PARSE TREE (Derivation Tree) (PT)

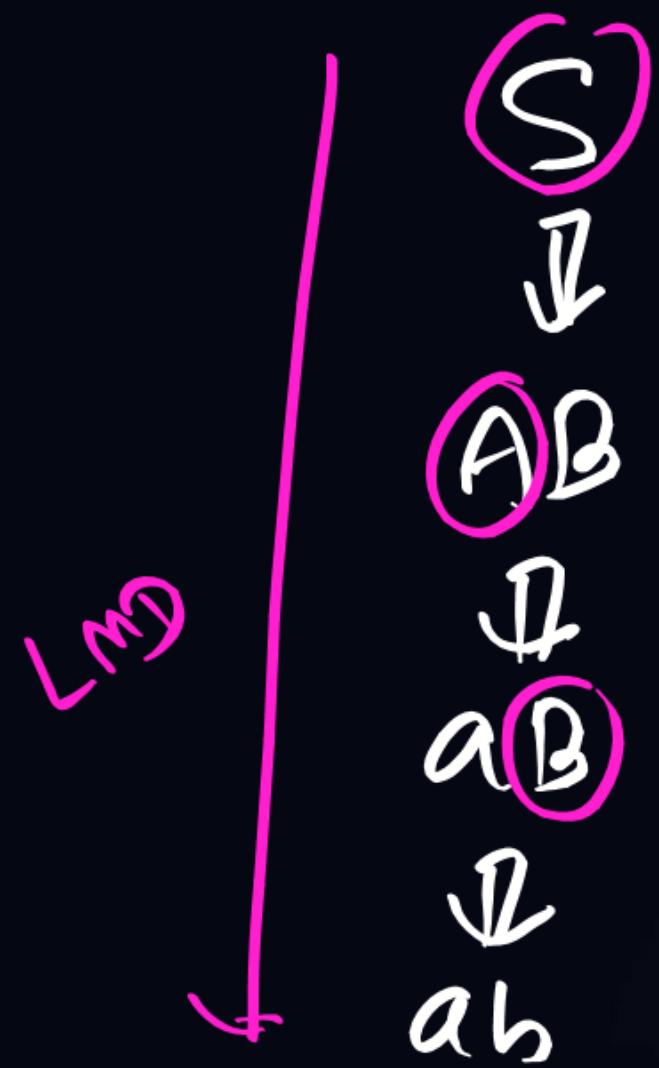
$S \rightarrow AB$   
 $A \rightarrow a$   
 $B \rightarrow b$

$w = ab$  starts

CFG



PW

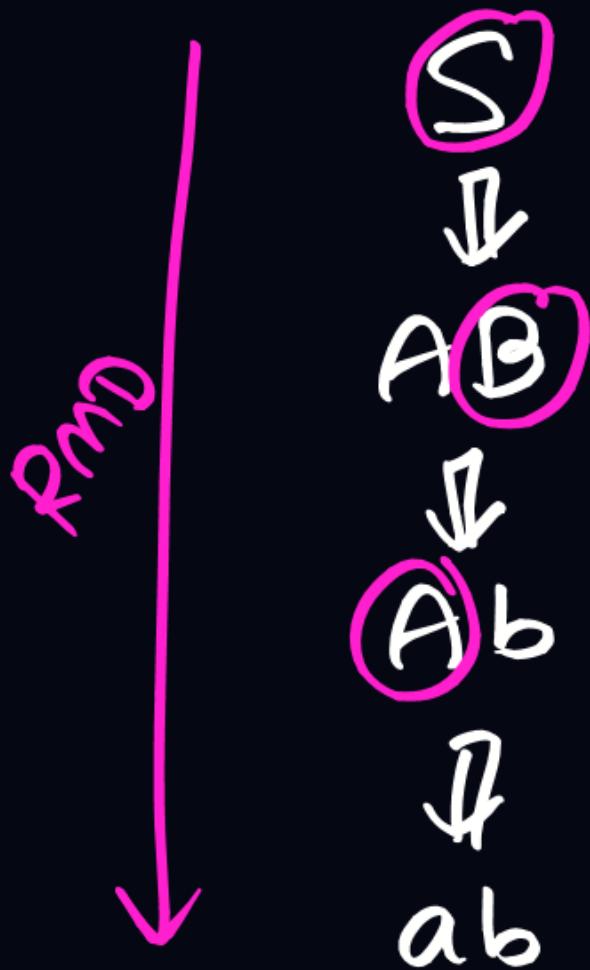


Length of Derivation  
= Derivation Length  
= No. of steps  
= No. of substitutions

= 3

What is LMD order?

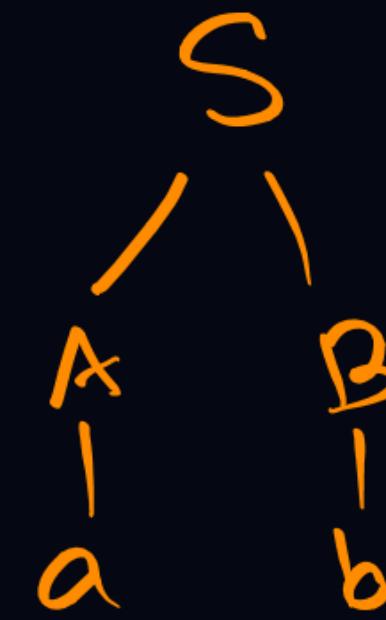
= SAB



Length of derivation = 3

RMD order :  
S, B, A

Note: LMD and RMD need not be same



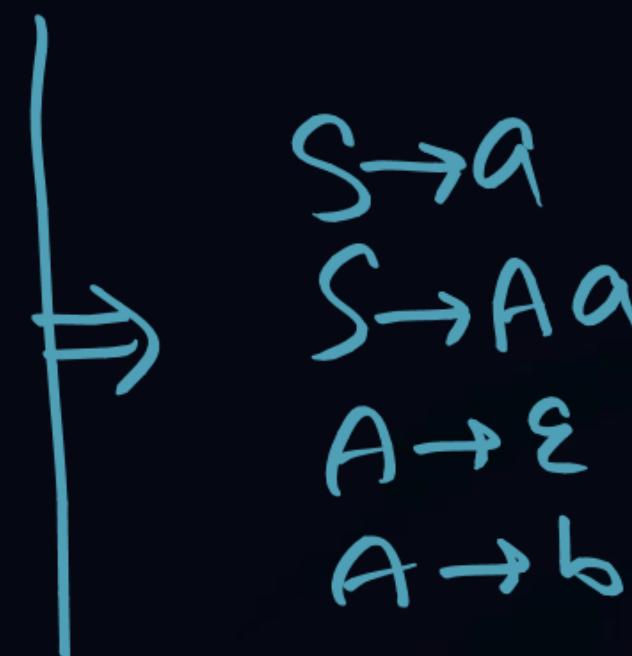
Derivation Length for "ab"

=  
No. of non leaf nodes

=

3

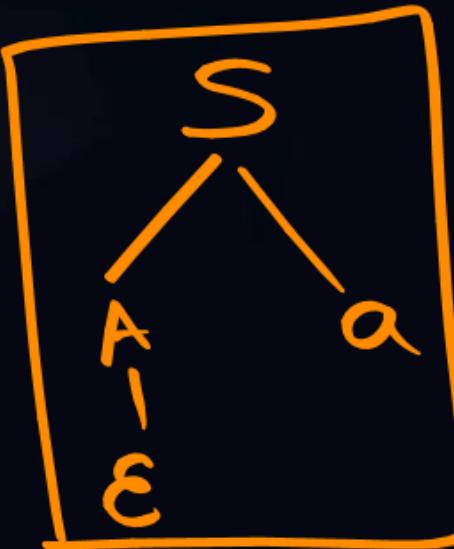
$$\begin{array}{l} S \rightarrow a \mid Aa \\ A \rightarrow \epsilon \mid b \end{array}$$



$w = a$  ⇒ How many derivations for  $w$  ?



OR

 $= a$  $\epsilon a = a$ 

= 2 derivations

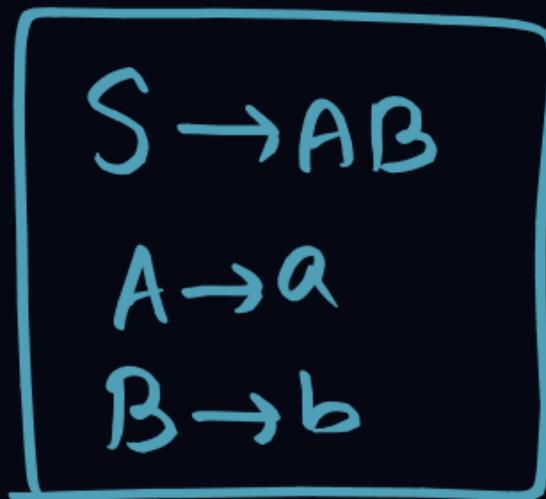
= 2 Parse Trees

= 2 LMDs

= 2 RMDs

For particular string:

$$\boxed{\begin{array}{l} \text{No. of derivations} \\ = \\ \text{No. of PTs} \\ = \\ \text{No. of LMDs} \\ = \\ \text{No. of RMDs} \end{array}}$$



$$w = ab$$

No. of LMDs for  $w = ab$

= No. of RMDs for  $w = ab$

$$= 1$$



In this example :

LMD & RMD are different

No. of steps in LMD

= No. of steps in RMD

$$= 3$$

**LMD**

**RMD**



①  $S \rightarrow SS | \epsilon$

P  
W

$W = \epsilon$

No. of derivations = ?

Infinite paths



=  $\epsilon$



$\epsilon \epsilon = \epsilon$



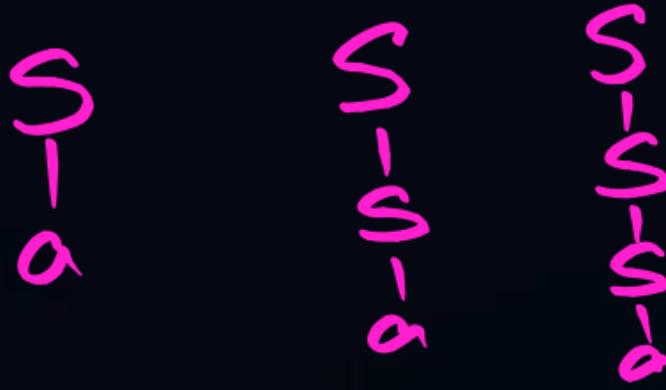
$\epsilon \epsilon \epsilon = \epsilon$

②

$$S \rightarrow S|a|ab$$

$$\begin{array}{l} S \rightarrow S \\ S \rightarrow a \\ S \rightarrow ab \end{array}$$

How many Parse Trees for "a" ? = Infinite



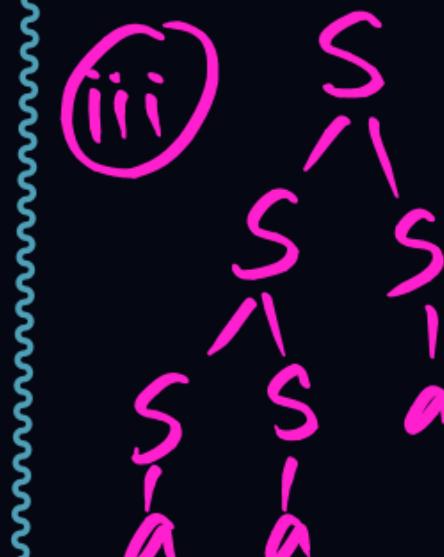
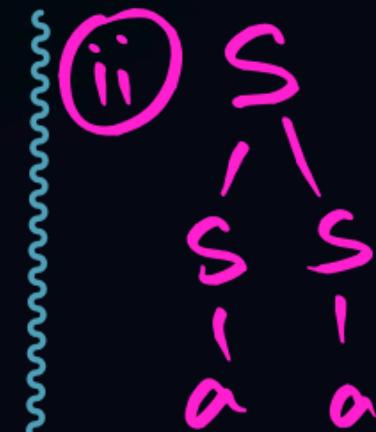
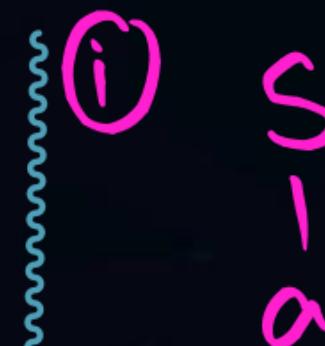
= Infinite

P  
W

③

$$S \rightarrow SS|a$$

- i) How many PTs for "a" ? = 1
- ii) How many PTs for "aa" ? = 1
- iii) How many PTs for "aaa" ? = 2



④  $S \rightarrow SSS | a$

P  
W



i) How many PTs for "a" ? = 1

ii) " " " for "aa" ? = 0



iii) " " " for "aaa" ? = 1

# Types of CFGs based on derivations:

## ① Ambiguous CFG

$$S \rightarrow \epsilon | a | b | A$$

$$A \rightarrow a$$

there is some string,  
which has more than  
one derivation.

## ② Unambiguous CFG

$$S \rightarrow \epsilon | a | b$$

$$\epsilon \Rightarrow 1 \text{ PT}$$

$$a \Rightarrow 1 \text{ PT}$$

$$b \Rightarrow 1 \text{ PT}$$

Every string has  
only 1 PT

# Types of CFGs based on derivations:

① Ambiguous CFG

Some string,  $> 1$  PT

$\exists w \in L(CFG) \Rightarrow > 1$  PT

② Unambiguous CFG

Every string,  $\text{only } 1$  PT

$\forall w \in L(CFG) \Rightarrow 1$  PT

# Types of CFGs based on derivations:

① Ambiguous CFG



Easy to write

② Unambiguous CFG



Better

# Identify Ambiguous and Unambiguous CFGs.



①  $S \rightarrow a$

Ambiguous

$a \Rightarrow IPT$

Unambiguous ✓

②  $S \rightarrow a | \epsilon$

Ambiguous

$\epsilon \Rightarrow I$   
 $a \Rightarrow I$

Unambiguous ✓

③  $S \rightarrow AB | \epsilon$

Ambiguous

$\epsilon \Rightarrow IPT$   
 $a \Rightarrow IPT$

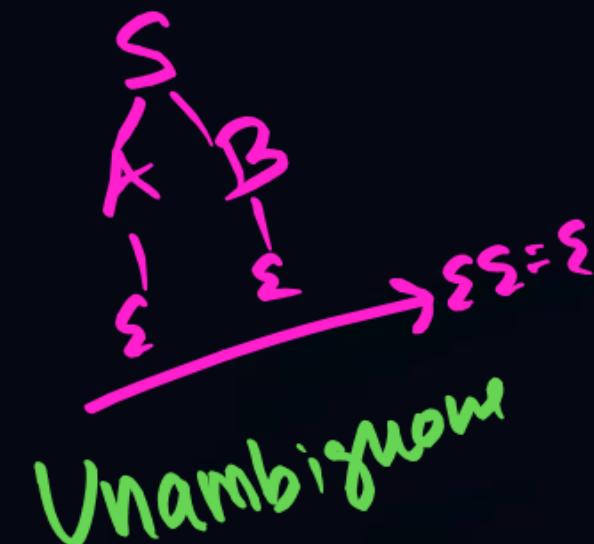
$A \rightarrow \epsilon$

$B \rightarrow a$

Unambiguous ✓

④  $S \rightarrow AB$   
 $A \rightarrow \epsilon$   
 $B \rightarrow \epsilon$

$\epsilon \Rightarrow 1 \text{ PT}$



⑦  $S \rightarrow AB$   
 $A \rightarrow aA | \epsilon$   
 $B \rightarrow aB | \epsilon$

Ambiguous

$\epsilon \Rightarrow 1 \text{ PT}$

$a \Rightarrow 2 \text{ PTs}$



⑧  $S \rightarrow SS | \epsilon$

Ambiguous

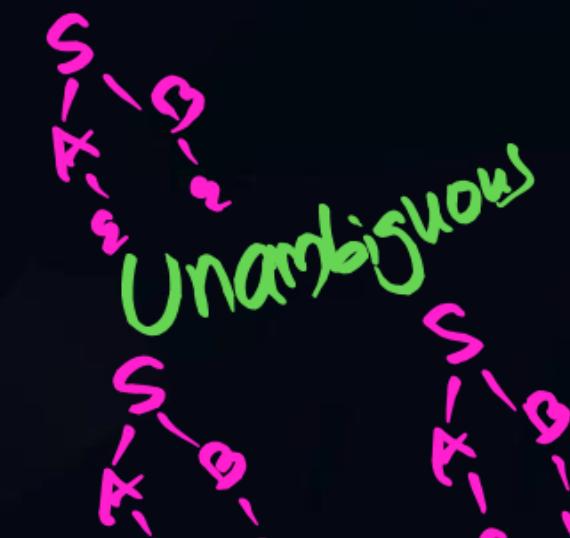
⑤  $S \rightarrow AB$   
 $A \rightarrow a | \epsilon$   
 $B \rightarrow b | \epsilon$

$\epsilon \Rightarrow 1 \text{ PT}$

$a \Rightarrow 1 \text{ PT}$

$b \Rightarrow 1 \text{ PT}$

~~$ab \Rightarrow 1 \text{ PT}$~~

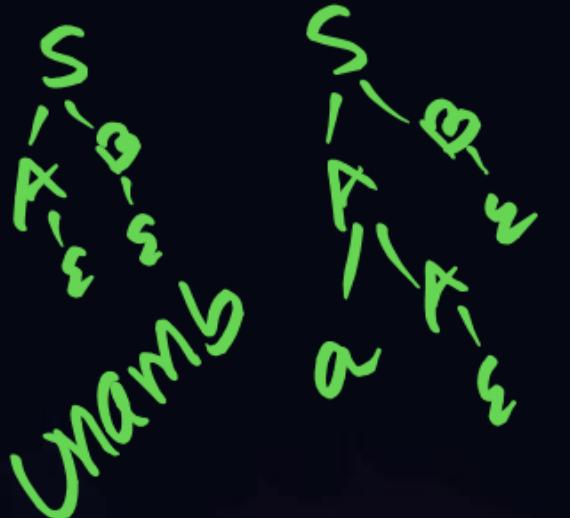


⑥  $S \rightarrow AB$   
 $A \rightarrow aA | \epsilon$   
 $B \rightarrow bB | \epsilon$

$\epsilon \Rightarrow 1 \text{ PT}$

$a \Rightarrow 1 \text{ PT}$

$b \Rightarrow 1 \text{ PT}$



⑨  $S \rightarrow SS | a$

Amb.

$a \Rightarrow 1 \text{ PT}$

$aa \Rightarrow 1 \text{ PT}$

$aaa$   $\Rightarrow 2 \text{ PTs}$

$S \rightarrow AB$  $A \rightarrow aA | \epsilon$  $B \rightarrow bB | \epsilon$ 

Unambiguous

 $\Sigma \rightarrow I$  $a \rightarrow I$  $b \rightarrow I$  $aa \rightarrow I$  $ab \rightarrow I$  $bb \rightarrow I$  $aaa \rightarrow I$  $aab \rightarrow I$  $abb \rightarrow I$  $bbb \rightarrow I$ 

I

I

⑩

$$S \rightarrow SaS | b$$

P  
W

⑪

$$S \rightarrow Sa | bS | c$$

⑫

$$E \rightarrow E + E | a$$

⑬

$$S \rightarrow AB | ab$$

$$A \rightarrow aA | \epsilon$$

$$B \rightarrow bB | \epsilon$$

H. W.

## Summary



CFG ✓

Derivations ✓

Amb CFG ✓

Unamb CFG ✓



THANK YOU GW  
SOLDIERS !