



## Chapter 2: High-Level Languages and Their Syntactic Description

Instructor: Gao Zhen

Email: [gaozhen@tongji.edu.cn](mailto:gaozhen@tongji.edu.cn)



# Outline

1. Programming language definition
2. General features of high-level languages
  - Program structure
  - Data types & operations
  - Statements & control structures
- 3. Syntax description methods**

# Formal Languages

- To understand/define high-level programming languages is crucial for compiler construction.
- This section will introduce the **formal description** of syntactic structures.

**Compiler = Formal language theory + Compiler techniques**

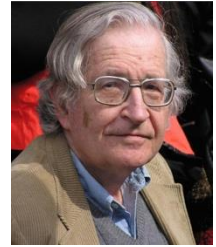
# Formal Languages

## ■ Formal languages

- Described using universally recognized **symbols** and expressions, formal languages are general and nationality-neutral.
- A formal language is a set of strings over some **alphabet** and has a **well-defined** descriptive scope.
- The original idea came from linguist **Noam Chomsky**, who aimed to use formal methods to describe languages.
- Started in natural language research, but found wide application in computer science, especially in theoretical computer science.

# Timeline

- 1936 — Alan Turing
  - Proposed the Turing Machine, formalizing the concept of “computability.”
- 1951–1956 — Stephen Kleene
  - Developed the concept of regular sets / regular expressions.  
Proved the equivalence regular expressions  $\leftrightarrow$  finite automata.
- 1956 — **Noam Chomsky**
  - Published *Three Models for the Description of Language*. Proposed the Chomsky hierarchy (regular, context-free, context-sensitive, unrestricted grammars).
- 1959 — Michael O. Rabin & Dana Stewart Scott
  - Introduced Nondeterministic Finite Automata (NFA) , deepening automata theory. Received the Turing Award in 1976.
- 1960–1961 — John Backus, Peter Naur, et al.
  - Used BNF (Backus–Naur Form) in programming language design, directly applying Chomsky’s context-free grammar ideas.
- 1961–1963 — John E. Hopcroft, Jeffrey D. Ullman etc.
  - Systematically proved the equivalence between formal languages and automata.
- 1969 — Hopcroft & Ullman
  - Published *Formal Languages and Their Relation to Automata*, a classic textbook that unified the equivalence of grammars, automata, and languages.



# Formal Definition of Grammar

- A grammar  $G$  is a 4-tuple  $G=(V_N, V_T, S, \mathfrak{E})$ 
  - $V_N$  : a finite set of non-terminal symbols
  - $V_T$  : a finite set of terminal symbols,  $V_N \cap V_T = \Phi$
  - $S$  : the start symbol, and  $S \in V_N$
  - $\mathfrak{E}$  : a finite set of production rules of the form  $P \rightarrow \alpha$   
 $P \in (V_N \cup V_T)^* V_N (V_N \cup V_T)^*$ ,  $\alpha \in (V_N \cup V_T)^*$

# Example of Grammar

■ Let  $G_1 = (\{N\}, \{0, 1\}, N, \{N \rightarrow 0N, N \rightarrow 1N, N \rightarrow 0, N \rightarrow 1\})$

$V_N = \{N\}$ , [non-terminal symbols]

$V_T = \{0, 1\}$ , [terminal symbols / alphabet]

$S = N$ , [start symbol]

$\mathcal{P} = \{N \rightarrow 0N, N \rightarrow 1N, N \rightarrow 0, N \rightarrow 1\}$  [productions]

# Explanation of Grammar

- Terminal symbols ( $V_T$ )

- basic symbols in the **alphabet**, **individual** tokens

- Non-terminal symbols ( $V_N$ )

- need further definition, represent **grammatical concept** such as "arithmetic expression", "boolean expression", "procedure", etc.  
→ Non-terminals represent **sets** (not individual tokens)

- Start symbol  $S$

- represents **the starting point** of the language, and all strings derived from it form the language
  - a **non-terminal** that appears on the left side of at least one production rule



# Explanation of Grammar

- **Production rules**  $\mathcal{E}$  : define the syntax structure; written as

- $P \rightarrow \alpha$  (or  $P ::= \alpha$ )

where " $\rightarrow$ " means "is **defined** as"

$P$ ,  $\alpha$  are symbol strings

Example :  $S \rightarrow 0S1$ ,  $0S \rightarrow 01$

- Productions with the same left-hand side ( $P \rightarrow \alpha_1$ ,  $P \rightarrow \alpha_2$ ,  $P \rightarrow \alpha_3$ , ...,  $P \rightarrow \alpha_n$ ) can be written compactly  $P \rightarrow \alpha_1 | \alpha_2 | \alpha_3 | \dots | \alpha_n$  Each  $\alpha_i$  is a **candidate** of  $P$ .

# Notation Conventions

- $V_N$ : **uppercase** letters A, B, C, S, etc.
- $V_T$ : **lowercase** letters, digits 0–9, operators +, −, etc.
- $\alpha$ 、 $\beta$ 、 $\gamma$ : **strings** of grammar symbols,  $\in (V_T \cup V_N)^*$
- $S$ : **start** symbol, appears in the first production
- $\rightarrow$  : **definition** symbol ("is defined as")
- $|$  : "or"

# Notation Conventions

- To simplify, only the production part is written  
Assume the left-hand side of the first production is the start symbol, or prefix the productions with "**G[A]**" where G is the grammar name, and A is the start symbol

Grammar G[N]:  $N \rightarrow 0N, N \rightarrow 1N, N \rightarrow 0, N \rightarrow 1$

Grammar G[E]:  $E \rightarrow E + E \mid E * E \mid (E) \mid i$

# Direct Derivation and Reduction

## ■ Direct derivation

- If  $A \rightarrow \gamma$  is a production, and  $\alpha, \beta \in (V_T \cup V_N)^*$ , then applying the rule  $A \rightarrow \gamma$  to string  $\alpha A \beta$  yields  $\alpha \gamma \beta$
- Written as:  $\alpha A \beta \Rightarrow \alpha \gamma \beta$ , called a **direct derivation**

## ■ Direct reduction

- Reduction is the reverse of derivation
- If  $\alpha A \beta \Rightarrow \alpha \gamma \beta$ , then  $\alpha \gamma \beta$  can be **directly reduced** to  $\alpha A \beta$

# Derivation

- Let  $\alpha_1, \alpha_2, \dots, \alpha_n$  ( $n > 0$ )  $\in (V_T \cup V_N)^*$ , And

$$\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$$

This sequence is a derivation from  $\alpha_1$  to  $\alpha_n$

- If such a derivation exists,  $\alpha_1$  can derive  $\alpha_n$

- $\alpha_1 \overset{+}{\Rightarrow} \alpha_n$  (via one or more steps)
- $\alpha_1 \overset{*}{\Rightarrow} \alpha_n$  (via zero or more steps)

# Multiple Derivations

- Given Grammar  $G[N_1]$ :

$N_1 \rightarrow N$      $N \rightarrow ND|D$      $D \rightarrow 0|1|2$

Then sentence "12" can be derived in multiple ways:

(1)  $N_1 \Rightarrow N \Rightarrow ND \Rightarrow N2 \Rightarrow D2 \Rightarrow 12$

(2)  $N_1 \Rightarrow N \Rightarrow ND \Rightarrow DD \Rightarrow 1D \Rightarrow 12$

- Thus, the same sentence can have different derivation sequences.

**Example – Grammar  $G[S]$ :**  $S \rightarrow AB$      $A \rightarrow A0|1B$   
 $B \rightarrow 0|S1$ ,    Generate the sentence **101001**:

**Leftmost derivation:**

$S \Rightarrow AB$   
 $\Rightarrow 1BB$   
 $\Rightarrow 10B$   
 $\Rightarrow 10S1$   
 $\Rightarrow 10AB1$   
 $\Rightarrow 101BB1$   
 $\Rightarrow 1010B1$   
 $\Rightarrow 101001$

**Rightmost derivation:**

$S \Rightarrow AB$   
 $\Rightarrow AS1$   
 $\Rightarrow AAB1$   
 $\Rightarrow AA01$   
 $\Rightarrow A1B01$   
 $\Rightarrow A1001$   
 $\Rightarrow 1B1001$   
 $\Rightarrow 101001$

# Sentential Forms, Sentences, and Language

$$G[E]: E \rightarrow E + E \mid E * E \mid (E) \mid i$$

- **Sentential Form:** Suppose  $G$  is a grammar and  $E$  is its starting symbol. If  $E \Rightarrow \alpha$ , then  $\alpha$  is a sentential form of grammar  $G$

Example:  $(E + E)$ ,  $(i + E)$ ,  $(i + i)$ ,  $E$

- **Sentence:** Sentential form consisting only of terminal characters are called sentences

Example:  $(i \times i + i)$ ,  $(i + i)$

- **Language:** The whole of the sentence produced by grammar  $G$

$$L(G) = \{\alpha \mid S \xRightarrow{+} \alpha, \alpha \in V_T^*\}$$



**Example.** With grammar G1 [S]:  $S \rightarrow bA$ ,  $A \rightarrow aA \mid a$   
Try to find the **language** described by this grammar.

**Solution:** Because the following sentences can be deduced from the start symbol S:

$S \Rightarrow bA \Rightarrow ba$

$S \Rightarrow bA \Rightarrow baA \Rightarrow baa$

$S \Rightarrow bA \Rightarrow baA \Rightarrow baaA \Rightarrow baaa$

...

$S \Rightarrow bA \Rightarrow baA \Rightarrow \dots \Rightarrow baa\dots a$

So  $L(G1) = \{ ba^n \mid n \geq 1 \}$

**Example.** Construct a grammar  $G_3$  for the language:  $L(G_3) = \{b^n a^n \mid n \geq 1\}$

This language consists of strings like:

$n=1, x=ba$

$n=2, x=bbaa$

$n=3, x=bbbbaaa \dots$

$L(G_3) = \{ba, bbaa, bbbbaaa, \dots\}$

Features of  $L(G_3)$ :

1. Symmetric strings (a's and b's in pairs)
2. Infinite set with recursive pattern
3. Alphabet:  $\{a, b\}$

**$G_3: S \rightarrow bSa \mid ba$**

**$G_3 = ( \{S\}, \{a, b\}, S, \{S \rightarrow bSa \mid ba\} )$**

# Exercise

## ■ P36

### □ 第6题

#### □ 令文法 $G_6$ 为

$N \rightarrow D | ND$

$D \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

(1)  $G_6$ 的语言 $L(G_6)$ 是什么？

(2) 给出句子0127, 34和568的最左推导和最右推导

### □ 第8题

#### □ 令文法为

$E \rightarrow T | E + T | E - T$

$T \rightarrow F | T * F | T / F$

$F \rightarrow (E) / i$

(1) 给出 $i+i*i$ ,  $i*(i+i)$ 的最左推导和最右推导

(2) 给出 $i+i+i$ ,  $i+i*i$ 和 $i-i-i$ 的语法树

# Chomsky Grammar System (Review)

- In the Chomsky hierarchy, any grammar must include:
  - Two distinct finite sets of symbols:
    - Non-terminal set  $V_N$
    - Terminal set  $V_T$
  - A start symbol  $S$
  - A finite set of formal rules  $\mathfrak{E}$  (productions)
- Grammar  $\mathbf{G}=(V_N, V_T, S, \mathfrak{E})$ ,  $\mathfrak{E} : P \rightarrow \alpha$ ,  
where  $P \in (V_N \cup V_T)^* V_N (V_N \cup V_T)^*$ ,  $\alpha \in (V_N \cup V_T)^*$
- Restrictions are imposed on the form of productions, and grammars are classified into four types: **Type 0**, **Type 1**, **Type 2**, and **Type 3**

# Type 0 Grammar

- Type 0 Grammar: Unrestricted grammar, phrase structure grammar
  - Corresponding Language: Recursively Enumerable Language  
Equivalent to Turing machines
- **Example: The following grammar is Type 0:**

$S \rightarrow aBC|aSBC$

$CB \rightarrow BC$

$aB \rightarrow ab$

$bB \rightarrow bb$

$bB \rightarrow b$

$bC \rightarrow bc$

$cC \rightarrow cc$

$cC \rightarrow c$

# Type 1 Grammar

- Also called Context-Sensitive Grammar (CSG)

- Production format:  $P \rightarrow \alpha$ , where

$$|P| \leq |\alpha|, \quad \alpha \in (V_N \cup V_T)^*, \quad P \in (V_N \cup V_T)^* V_N (V_N \cup V_T)^*$$

- Corresponding language: **CSL** (Context-sensitive Language)
- If  $\varepsilon$  is not considered, it is equivalent to a Linear Bounded Automata (**LBA**)
- Example :  $\alpha A \beta \rightarrow \alpha \gamma \beta$

# Type 1 Grammar

- **Example:** The following grammar is Type 1:

$S \rightarrow aBC | aSBC$

$CB \rightarrow BC$

$aB \rightarrow ab$

$bB \rightarrow bd$

$bC \rightarrow bc$

$cC \rightarrow ce$

# Type 2 Grammar

- Also called Context-free Grammar (CFG)

- Production form:  $P \rightarrow \alpha$ , where

$$P \in V_N, \quad \alpha \in (V_N \cup V_T)^*$$

- Corresponding language: Context-free Language (CFL)
- Corresponding automaton: Pushdown Automata (PDA)



# Type 2 Grammar

## ■ Example. Context-free Grammar

$S \rightarrow 01$

$S \rightarrow 0S1$

Generated language  $L = \{ 0^n 1^n \mid n \geq 1 \}$ ,

e.g.  $0011, 000111, 01 \in L$ , but  $10, 1001, \varepsilon, 010 \notin L$ .

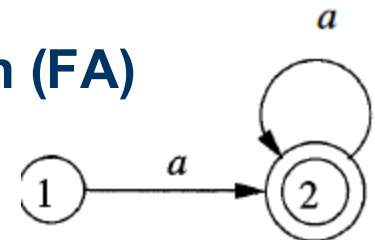
No finite automata can accept  $L$ .

# Type 3 Grammar

## ■ Also called Regular Grammar

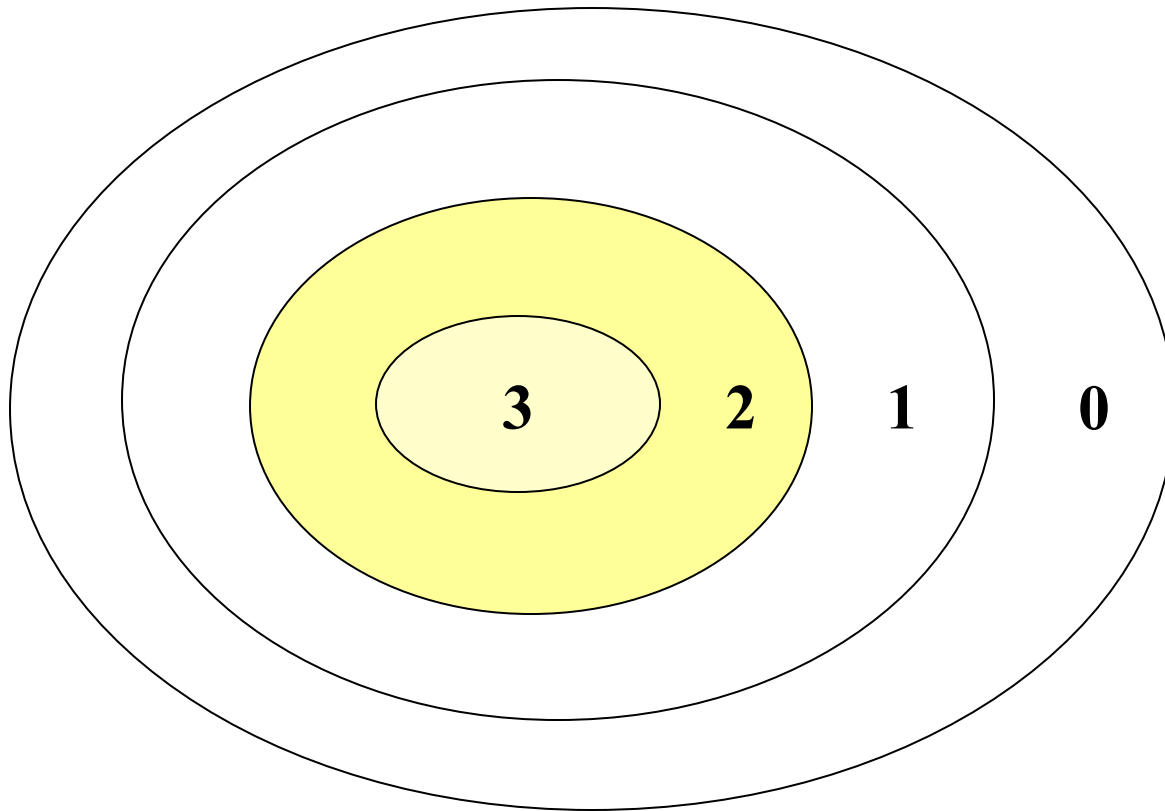
- Right-linear Grammar : Productions of the form  $A \rightarrow \omega B$  or  $A \rightarrow \omega$ , where  $A, B \in V_N$ ,  $\omega \in V_T^*$ .
- Left-linear Grammar: Productions of the form  $A \rightarrow B\omega$  or  $A \rightarrow \omega$ , where  $A, B \in V_N$ ,  $\omega \in V_T^*$ .
- Equivalent to Regular Expressions
- Corresponding language: Regular Language
- Corresponding automaton: Finite Automaton (FA)

## ■ Example $S \rightarrow aS, S \rightarrow a$

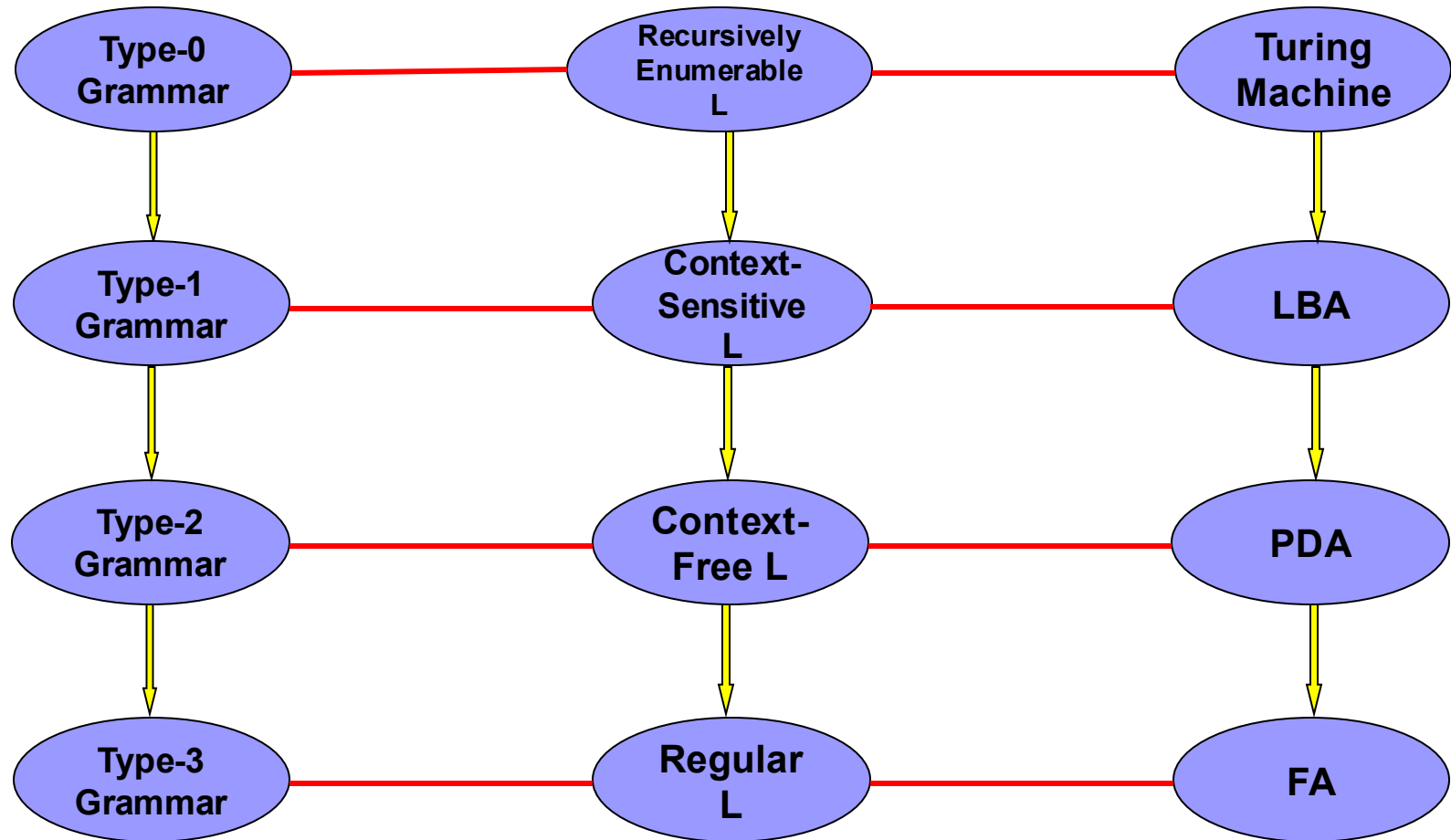


Equivalent regular expression :  $a^+$ , or  $a^*a$

# Chomsky Grammar Hierarchy



# Correspondence Between Grammars, Formal Languages, and Automata



## Exercise : how to draw the Finite Automata

**Example. Given grammar**

**G1[S]:**  $S \rightarrow bA$  ,  $A \rightarrow aA \mid a$

**Task:** Draw the equivalent Finite Automata

# Conclusion (1)

- **Restrictions** on context-free grammars for programming languages
  - (1) Productions that form “**self-loops**” are **NOT** allowed:  $P \rightarrow P$
  - (2) Every non-terminal  $P$  is **useful**,  $P \in V_N$  :
    - Appears in some derivation from start symbol:  $S \Rightarrow \alpha P \beta$
    - Can derive a terminal string:  $P \Rightarrow \gamma, \gamma \in V_T^*$
- **Lexical analysis: based on regular grammars**

**Regular Grammar**

$$A \rightarrow IB \qquad B \rightarrow IB|dB|\epsilon$$
- **Syntax analysis: based on context-free grammars**

# Conclusion (2)

- Context-free grammars are powerful enough to describe the syntax of most modern programming languages

- Arithmetic Expression
- Assignment Statement
- Conditional Statement
- .....

## Arithmetic Expression

Grammar  $G = (\{E\}, \{+, *, i, (, )\}, E, P)$

$E \rightarrow i$                        $E \rightarrow E + E$

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

## Conditional Statement

$S \rightarrow \text{if } E \text{ then } S$

$S \rightarrow \text{if } E \text{ then } S \text{ else } S$

# Syntax Tree and Grammar Ambiguity

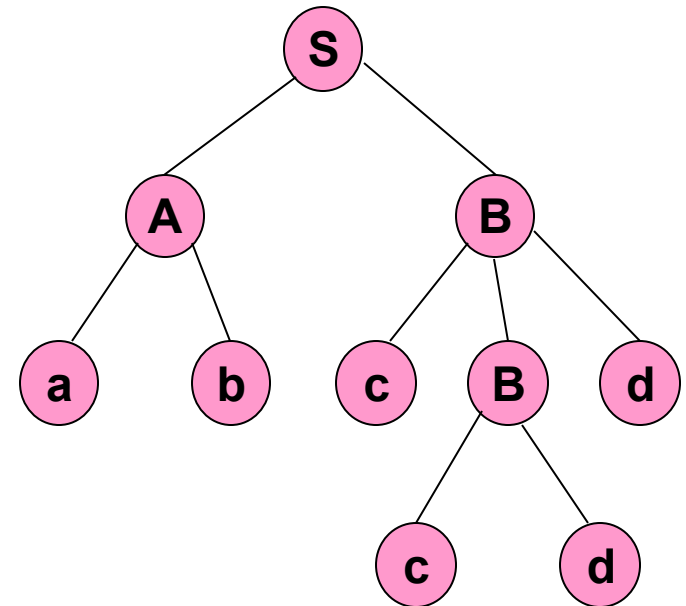
- **Parse Tree (Syntax Tree)**
- **Ambiguity**





# What is a Parse Tree?

- A parse tree represents the derivation of a sentence (string) in a grammar
- It is an inverted tree (root at the top, leaves at the bottom)
  - Node
  - Edge
  - Root Node
  - Leaf Node
  - Leaf Branches :  $A(a,b)$ ,  $B(c,d)$
  - Sibling Nodes



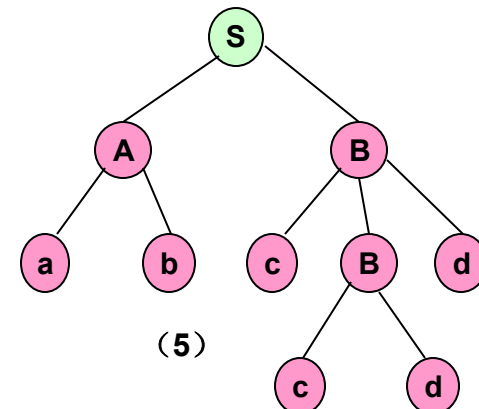
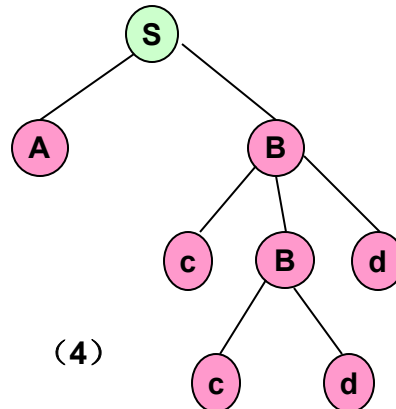
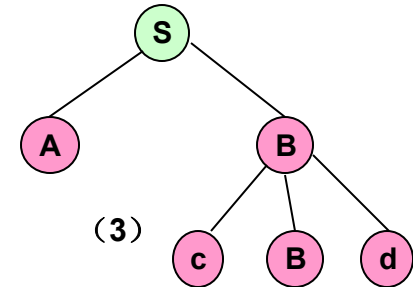
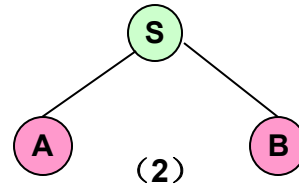
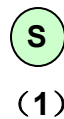
# Constructing a Parse Tree from a Derivation

$S \Rightarrow AB$

$\Rightarrow AcBd$

$\Rightarrow Accdd$

$\Rightarrow abccdd$



■ Example. Grammar  $G[E]$  :

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid i$$

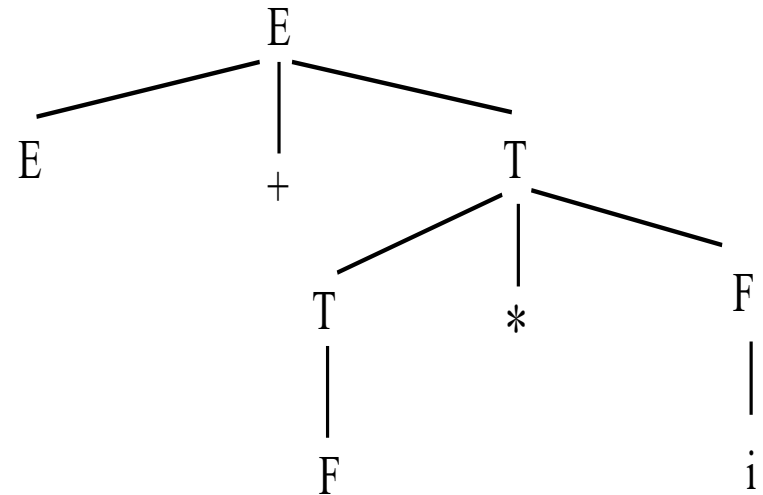
Derivation of expression  $E + F * i$  :

$$E \Rightarrow E + T$$

$$\Rightarrow E + T * F$$

$$\Rightarrow E + F * F$$

$$\Rightarrow E + F * i$$



The syntax tree

# More about parse tree

- A parse tree shows which rules are applied and on which nonterminal symbols, but it does **NOT** indicate **the order of rule applications**
- Does a sentence have a **unique** leftmost (or rightmost) derivation?
  - Not always
- Does a sentence correspond to a **unique** parse tree?
  - Not always

# Ambiguity

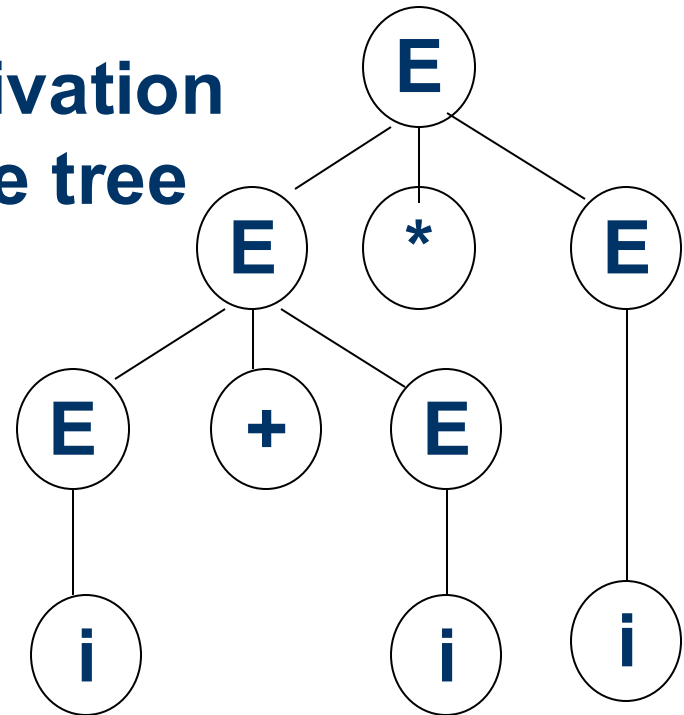
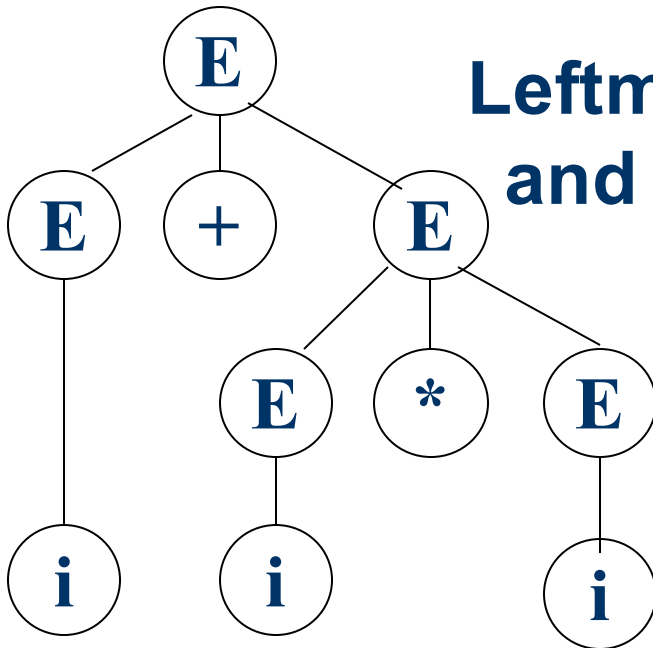
- A sentence is ambiguous if it has **two** parse trees
- A grammar is ambiguous if it generates **any** ambiguous sentence; otherwise, it is unambiguous

Grammar  $G[E] : E \rightarrow E + E \mid E * E \mid (E) \mid i$

$E \Rightarrow E + E \Rightarrow i + E$   
 $\Rightarrow i + E * E \Rightarrow i + i * E$   
 $\Rightarrow i + i * i$

$E \Rightarrow E * E \Rightarrow E + E * E$   
 $\Rightarrow i + E * E \Rightarrow i + i * E$   
 $\Rightarrow i + i * i$

Leftmost derivation  
and its parse tree



Similarly, the sentence's rightmost derivation and its parse tree are also different.

# Transforming an ambiguous grammar into an unambiguous grammar

- Grammar  $G[E] : E \rightarrow E + E \mid E \times E \mid (E) \mid i$  is ambiguous. By defining precedence ( $\times > +$ ) and left associativity, it can be transformed into an unambiguous grammar.

Grammar  $G[E]$ :

$$E \rightarrow T \mid E + T$$
$$T \rightarrow F \mid T \times F$$
$$F \rightarrow ( E ) \mid i$$

The sentence has a unique derivation :  $(i \times i + i)$



# QUIZ-CANVAS



**Dank u**

Dutch

**Merci**

French

**Спасибо**

Russian

**Gracias**

Spanish

شكراً

Arabic

감사합니다

Korean

תודה רבה

Hebrew

**Tack så mycket**

Swedish

धन्यवाद

Hindi

**Obrigado**

Brazilian  
Portuguese

谢谢

Chinese

***Thank You !***

**Dankon**

Esperanto

ありがとうございます

Japanese

**Trugarez**

Breton

**Danke**

German

**Tak**

Danish

**Grazie**

Italian

நன்றி

Tamil

**děkuji**

Czech

ขอบคุณ

Thai

**go raibh maith agat**

Gaelic