

# Context-Free Grammars

**Example:** We would like to describe a language of arithmetic expressions, containing expressions such as:

175      (9+15)      (((10-4)\*((1+34)+2))/(3+(-37)))

For simplicity we assume that:

- Expressions are fully parenthesized.
- The only arithmetic operations are “+”, “-”, “\*”, “/” and unary “-”.
- Values of operands are natural numbers written in decimal — a number is represented as a non-empty sequence of digits.

Alphabet:  $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, -, *, /, (, )\}$

**Example (cont.):** A description by an inductive definition:

- **Digit** is any of characters 0, 1, 2, 3, 4, 5, 6, 7, 8, 9.
- **Number** is a non-empty sequence of digits, i.e.:
  - If  $\alpha$  is a digit then  $\alpha$  is a number.
  - If  $\alpha$  is a digit and  $\beta$  is a number then also  $\alpha\beta$  is a number.
- **Expression** is a sequence of symbols constructed according to the following rules:
  - If  $\alpha$  is a number then  $\alpha$  is an expression.
  - If  $\alpha$  is an expression then also  $(-\alpha)$  is an expression.
  - If  $\alpha$  and  $\beta$  are expressions then also  $(\alpha+\beta)$  is an expression.
  - If  $\alpha$  and  $\beta$  are expressions then also  $(\alpha-\beta)$  is an expression.
  - If  $\alpha$  and  $\beta$  are expressions then also  $(\alpha*\beta)$  is an expression.
  - If  $\alpha$  and  $\beta$  are expressions then also  $(\alpha/\beta)$  is an expression.

# Context-Free Grammars

**Example (cont.):** The same information that was described by the previous inductive definition can be represented by a **context-free grammar**:

New auxiliary symbols, called **nonterminals**, are introduced:

- $D$  — stands for an arbitrary digit
- $C$  — stands for an arbitrary number
- $E$  — stands for an arbitrary expression

$$D \rightarrow 0$$

$$D \rightarrow 1$$

$$D \rightarrow 2$$

$$D \rightarrow 3$$

$$D \rightarrow 4$$

$$D \rightarrow 5$$

$$D \rightarrow 6$$

$$D \rightarrow 7$$

$$D \rightarrow 8$$

$$D \rightarrow 9$$

$$C \rightarrow D$$

$$C \rightarrow DC$$

$$E \rightarrow C$$

$$E \rightarrow (-E)$$

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

$$E \rightarrow (E-E)$$

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

$$E \rightarrow (E / E)$$

**Example (cont.):** Written in a more succinct way:

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

$$C \rightarrow D \mid DC$$

$$E \rightarrow C \mid (-E) \mid (E+E) \mid (E-E) \mid (E * E) \mid (E/E)$$

**Example:** A language where words are (possibly empty) sequences of expressions described in the previous example, where individual expressions are separated by commas (the alphabet must be extended with symbol “,”):

$$S \rightarrow T \mid \varepsilon$$

$$T \rightarrow E \mid E, T$$

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

$$C \rightarrow D \mid DC$$

$$E \rightarrow C \mid (-E) \mid (E+E) \mid (E-E) \mid (E * E) \mid (E/E)$$

# Context-Free Grammars

**Example:** Statements of some programming language (a fragment of a grammar):

$$\begin{aligned} S &\rightarrow E; \mid T \mid \text{if } (E) \ S \mid \text{if } (E) \ S \ \text{else } S \\ &\quad \mid \text{while } (E) \ S \mid \text{do } S \ \text{while } (E); \mid \text{for } (F; F; F) \ S \\ &\quad \mid \text{return } F; \\ T &\rightarrow \{ \ U \} \\ U &\rightarrow \varepsilon \mid SU \\ F &\rightarrow \varepsilon \mid E \\ E &\rightarrow \dots \end{aligned}$$

**Remark:**

- $S$  — statement
- $T$  — block of statements
- $U$  — sequence of statements
- $E$  — expression
- $F$  — optional expression that can be omitted

# Context-Free Grammars

Formally, a **context-free grammar** is a tuple

$$\mathcal{G} = (\Pi, \Sigma, S, P)$$

where:

- $\Pi$  is a finite set of **nonterminal symbols** (**nonterminals**)
- $\Sigma$  is a finite set of **terminal symbols** (**terminals**),  
where  $\Pi \cap \Sigma = \emptyset$
- $S \in \Pi$  is an **initial nonterminal**
- $P \subseteq \Pi \times (\Pi \cup \Sigma)^*$  is a finite set of **rewrite rules**



## Remarks:

- We will use uppercase letters  $A, B, C, \dots$  to denote nonterminal symbols.
- We will use lowercase letters  $a, b, c, \dots$  or digits  $0, 1, 2, \dots$  to denote terminal symbols.
- We will use lowercase Greek letters  $\alpha, \beta, \gamma, \dots$  to denote strings from  $(\Pi \cup \Sigma)^*$ .
- We will use the following notation for rules instead of  $(A, \alpha)$

$$A \rightarrow \alpha$$

$A$  – left-hand side of the rule

$\alpha$  – right-hand side of the rule

**Example:** Grammar  $\mathcal{G} = (\Pi, \Sigma, S, P)$  where

- $\Pi = \{A, B, C\}$
- $\Sigma = \{a, b\}$
- $S = A$
- $P$  contains rules

$$A \rightarrow aBBb$$

$$A \rightarrow AaA$$

$$B \rightarrow \varepsilon$$

$$B \rightarrow bCA$$

$$C \rightarrow AB$$

$$C \rightarrow a$$

$$C \rightarrow b$$

**Remark:** If we have more rules with the same left-hand side, as for example

$$A \rightarrow \alpha_1$$

$$A \rightarrow \alpha_2$$

$$A \rightarrow \alpha_3$$

we can write them in a more succinct way as

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \alpha_3$$

For example, the rules of the grammar from the previous slide can be written as

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

*A*

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$\underline{A} \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

A

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$\underline{A} \rightarrow \underline{aBBb} \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$\underline{A} \Rightarrow \underline{aBBb}$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb$$



# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$\underline{B} \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow a\underline{B}Bb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$\underline{B} \rightarrow \varepsilon \mid \underline{bCA}$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow a\underline{B}Bb \Rightarrow a\underline{bCA}Bb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$\underline{A} \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abC\underline{A}Bb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$\underline{A} \rightarrow \underline{aBBb} \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abC\underline{ABb} \Rightarrow abC\underline{aBBb}Bb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$\underline{B} \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaB\underline{B}bBb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$\underline{B} \rightarrow \underline{\varepsilon} \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaB\underline{B}bBb \Rightarrow abCaBbBb$$



# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$\underline{C} \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow ab\underline{C}aBbBb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$\underline{C} \rightarrow AB \mid a \mid \underline{b}$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow ab\underline{C}aBbBb \Rightarrow ab\underline{b}aBbBb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$\underline{B} \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBb\underline{B}b$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$\underline{B} \rightarrow \underline{\varepsilon} \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBb\underline{B}b \Rightarrow abbaBbb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb \Rightarrow abbaBbb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$\underline{B} \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb \Rightarrow abba\underline{B}bb$$



# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$\underline{B} \rightarrow \underline{\varepsilon} \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb \Rightarrow abba\underline{B}bb \Rightarrow abbabb$$

# Context-Free Grammars

Grammars are used for generating words.

**Example:**  $\mathcal{G} = (\Pi, \Sigma, A, P)$  where  $\Pi = \{A, B, C\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains rules

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

For example, the word *abbabb* can be in grammar  $\mathcal{G}$  generated as follows:

$$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb \Rightarrow abbaBbb \Rightarrow abbabb$$

# Context-Free Grammars

On strings from  $(\Pi \cup \Sigma)^*$  we define relation  $\Rightarrow \subseteq (\Pi \cup \Sigma)^* \times (\Pi \cup \Sigma)^*$  such that

$$\alpha \Rightarrow \alpha'$$

iff  $\alpha = \beta_1 A \beta_2$  and  $\alpha' = \beta_1 \gamma \beta_2$  for some  $\beta_1, \beta_2, \gamma \in (\Pi \cup \Sigma)^*$  and  $A \in \Pi$  where  $(A \rightarrow \gamma) \in P$ .

**Example:** If  $(B \rightarrow bCA) \in P$  then

$$aCBbA \Rightarrow aCbCAbA$$

**Remark:** Informally,  $\alpha \Rightarrow \alpha'$  means that it is possible to derive  $\alpha'$  from  $\alpha$  by one step where an occurrence of some nonterminal  $A$  in  $\alpha$  is replaced with the right-hand side of some rule  $A \rightarrow \gamma$  with  $A$  on the left-hand side.

# Context-Free Grammars

On strings from  $(\Pi \cup \Sigma)^*$  we define relation  $\Rightarrow \subseteq (\Pi \cup \Sigma)^* \times (\Pi \cup \Sigma)^*$  such that

$$\alpha \Rightarrow \alpha'$$

iff  $\alpha = \beta_1 A \beta_2$  and  $\alpha' = \beta_1 \gamma \beta_2$  for some  $\beta_1, \beta_2, \gamma \in (\Pi \cup \Sigma)^*$  and  $A \in \Pi$  where  $(A \rightarrow \gamma) \in P$ .

**Example:** If  $(B \rightarrow bCA) \in P$  then

$$aC\underline{B}bA \Rightarrow aC\underline{bCA}bA$$

**Remark:** Informally,  $\alpha \Rightarrow \alpha'$  means that it is possible to derive  $\alpha'$  from  $\alpha$  by one step where an occurrence of some nonterminal  $A$  in  $\alpha$  is replaced with the right-hand side of some rule  $A \rightarrow \gamma$  with  $A$  on the left-hand side.

# Context-Free Grammars

A **derivation** of length  $n$  is a sequence  $\beta_0, \beta_1, \beta_2, \dots, \beta_n$ , where  $\beta_i \in (\Pi \cup \Sigma)^*$ , and where  $\beta_{i-1} \Rightarrow \beta_i$  for all  $1 \leq i \leq n$ , which can be written more succinctly as

$$\beta_0 \Rightarrow \beta_1 \Rightarrow \beta_2 \Rightarrow \dots \Rightarrow \beta_{n-1} \Rightarrow \beta_n$$

The fact that for given  $\alpha, \alpha' \in (\Pi \cup \Sigma)^*$  and  $n \in \mathbb{N}$  there exists some derivation  $\beta_0 \Rightarrow \beta_1 \Rightarrow \beta_2 \Rightarrow \dots \Rightarrow \beta_{n-1} \Rightarrow \beta_n$ , where  $\alpha = \beta_0$  and  $\alpha' = \beta_n$ , is denoted

$$\alpha \Rightarrow^n \alpha'$$

The fact that  $\alpha \Rightarrow^n \alpha'$  for some  $n \geq 0$ , is denoted

$$\alpha \Rightarrow^* \alpha'$$

**Remark:** Relation  $\Rightarrow^*$  is the reflexive and transitive closure of relation  $\Rightarrow$  (i.e., the smallest reflexive and transitive relation containing relation  $\Rightarrow$ ).

**Sentential forms** are those  $\alpha \in (\Pi \cup \Sigma)^*$ , for which

$$S \Rightarrow^* \alpha$$

where  $S$  is the initial nonterminal.

A **language**  $\mathcal{L}(\mathcal{G})$  generated by a grammar  $\mathcal{G} = (\Pi, \Sigma, S, P)$  is the set of all words over alphabet  $\Sigma$  that can be derived by some derivation from the initial nonterminal  $S$  using rules from  $P$ , i.e.,

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

**Example:** We want to construct a grammar generating the language

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



**Example:** We want to construct a grammar generating the language

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

Grammar  $\mathcal{G} = (\Pi, \Sigma, S, P)$  where  $\Pi = \{S\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains

$$S \rightarrow \varepsilon \mid aSb$$

# Context-Free Grammars

**Example:** We want to construct a grammar generating the language

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

Grammar  $\mathcal{G} = (\Pi, \Sigma, S, P)$  where  $\Pi = \{S\}$ ,  $\Sigma = \{a, b\}$ , and  $P$  contains

$$S \rightarrow \varepsilon \mid aSb$$

$$S \Rightarrow \varepsilon$$

$$S \Rightarrow aSb \Rightarrow ab$$

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aabb$$

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aaaSbbb \Rightarrow aaabbb$$

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aaaSbbb \Rightarrow aaaaSbbbb \Rightarrow aaaabbbb$$

...

**Example:** We want to construct a grammar generating the language consisting of all palindroms over the alphabet  $\{a, b\}$ , i.e.,

$$L = \{w \in \{a, b\}^* \mid w = w^R\}$$

**Remark:**  $w^R$  denotes the **reverse** of a word  $w$ , i.e., the word  $w$  written backwards.

**Example:** We want to construct a grammar generating the language consisting of all palindroms over the alphabet  $\{a, b\}$ , i.e.,

$$L = \{w \in \{a, b\}^* \mid w = w^R\}$$

**Remark:**  $w^R$  denotes the **reverse** of a word  $w$ , i.e., the word  $w$  written backwards.

*Solution:*

$$S \rightarrow \varepsilon \mid a \mid b \mid aSa \mid bSb$$

**Example:** We want to construct a grammar generating the language consisting of all palindroms over the alphabet  $\{a, b\}$ , i.e.,

$$L = \{w \in \{a, b\}^* \mid w = w^R\}$$

**Remark:**  $w^R$  denotes the **reverse** of a word  $w$ , i.e., the word  $w$  written backwards.

*Solution:*

$$S \rightarrow \varepsilon \mid a \mid b \mid aSa \mid bSb$$

$$S \Rightarrow aSa \Rightarrow abSba \Rightarrow abaSaba \Rightarrow abaaaba$$

**Example:** We want to construct a grammar generating the language  $L$  consisting of all correctly parenthesised sequences of symbols '(' and ')'. For example  $((())()) \in L$  but  $)() \notin L$ .

**Example:** We want to construct a grammar generating the language  $L$  consisting of all correctly parenthesised sequences of symbols '(' and ')'.

For example  $((())()) \in L$  but  $)() \notin L$ .

*Solution:*

$$S \rightarrow \varepsilon \mid (S) \mid SS$$

**Example:** We want to construct a grammar generating the language  $L$  consisting of all correctly parenthesised sequences of symbols '(' and ')'. For example  $((())()) \in L$  but  $)() \notin L$ .

*Solution:*

$$S \rightarrow \varepsilon \mid (S) \mid SS$$

$$\begin{aligned} S &\Rightarrow SS \Rightarrow (S)S \Rightarrow (S)(S) \Rightarrow (SS)(S) \Rightarrow ((S)S)(S) \Rightarrow \\ &((()S)(S)) \Rightarrow ((()S))S \Rightarrow ((()())S) \Rightarrow ((()())((S))) \Rightarrow \\ &((()())()) \end{aligned}$$



**Example:** We want to construct a grammar generating the language  $L$  consisting of all correctly constructed arithmetic expressions where operands are always of the form ' $a$ ' and where symbols  $+$  and  $*$  can be used as operators.

For example  $(a + a) * a + (a * a) \in L$ .

**Example:** We want to construct a grammar generating the language  $L$  consisting of all correctly constructed arithmetic expressions where operands are always of the form ' $a$ ' and where symbols  $+$  and  $*$  can be used as operators.

For example  $(a + a) * a + (a * a) \in L$ .

*Solution:*

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

# Context-Free Grammars

**Example:** We want to construct a grammar generating the language  $L$  consisting of all correctly constructed arithmetic expressions where operands are always of the form ' $a$ ' and where symbols  $+$  and  $*$  can be used as operators.

For example  $(a + a) * a + (a * a) \in L$ .

*Solution:*

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

$$\begin{aligned} E &\Rightarrow E + E \Rightarrow E * E + E \Rightarrow (E) * E + E \Rightarrow (E + E) * E + E \Rightarrow \\ &(a + E) * E + E \Rightarrow (a + a) * E + E \Rightarrow (a + a) * a + E \Rightarrow (a + a) * a + (E) \Rightarrow \\ &(a + a) * a + (E * E) \Rightarrow (a + a) * a + (a * E) \Rightarrow (a + a) * a + (a * a) \end{aligned}$$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$

# Derivation Tree

$A$

$$A \rightarrow aBBb \mid AaA$$

$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

$A$

# Derivation Tree

A

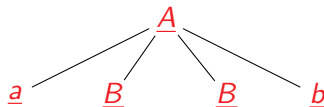
A  $\rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$

A

# Derivation Tree



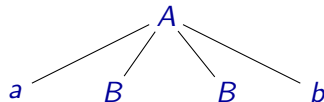
$\underline{A} \rightarrow \underline{aBBb} \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$

$\underline{A} \Rightarrow \underline{aBBb}$

# Derivation Tree



$$A \rightarrow aBBb \mid AaA$$

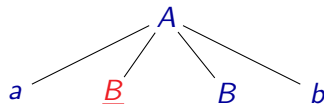
$$B \rightarrow \varepsilon \mid bCA$$

$$C \rightarrow AB \mid a \mid b$$

$$A \Rightarrow aBBb$$



# Derivation Tree



$A \rightarrow aBBb \mid AaA$

$\underline{B} \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$

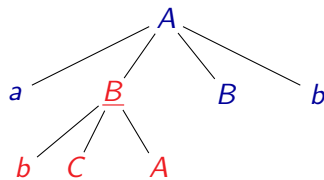
$A \Rightarrow a\underline{B}Bb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$\underline{B} \rightarrow \varepsilon \mid \underline{bCA}$

$C \rightarrow AB \mid a \mid b$



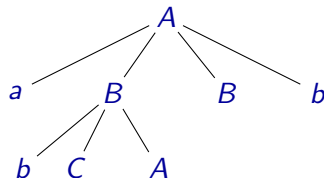
$A \Rightarrow a\underline{B}Bb \Rightarrow ab\underline{C}ABb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



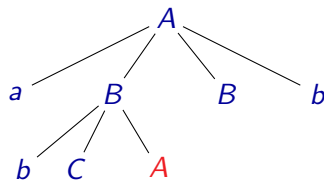
$A \Rightarrow aBBb \Rightarrow abCABb$

# Derivation Tree

A  $\rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



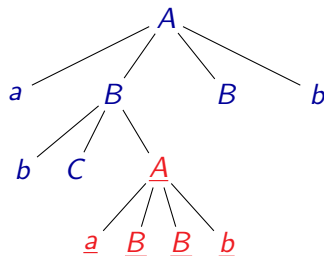
$A \Rightarrow aBBb \Rightarrow abC\underline{A}Bb$

# Derivation Tree

$\underline{A} \rightarrow \underline{aBBb} \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



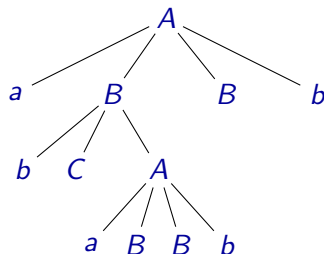
$A \Rightarrow aBBb \Rightarrow abC\underline{A}Bb \Rightarrow abC\underline{aBBb}Bb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



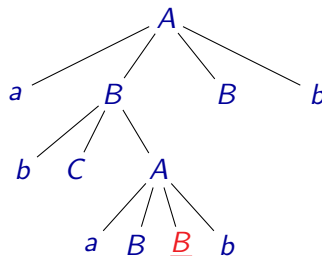
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$\underline{B} \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



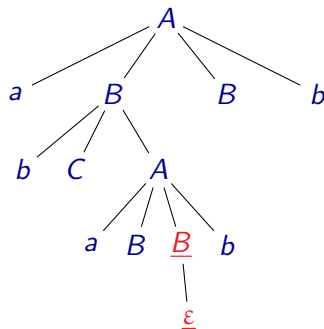
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaB\underline{B}bBb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$\underline{B} \rightarrow \underline{\epsilon} \mid bCA$

$C \rightarrow AB \mid a \mid b$



$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaB\underline{B}bBb \Rightarrow abCaBbBb$

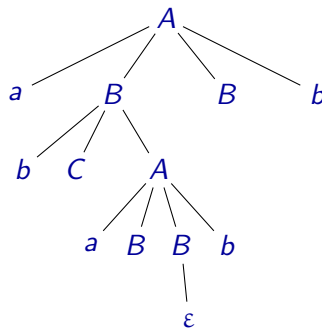


# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



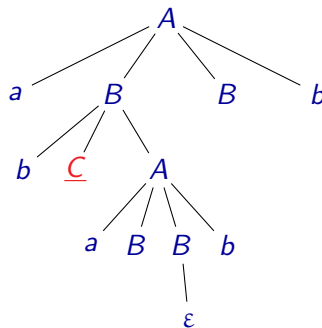
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C$   $\rightarrow AB \mid a \mid b$



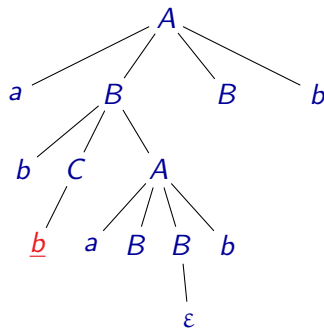
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow ab\underline{C}aBbBb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$\underline{C} \rightarrow AB \mid a \mid \underline{b}$



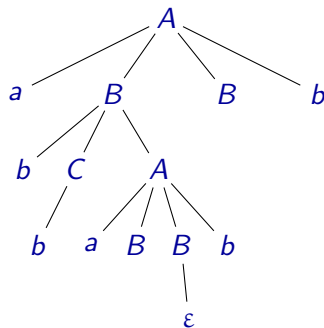
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow ab\underline{C}aBbBb \Rightarrow ab\underline{b}aBbBb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



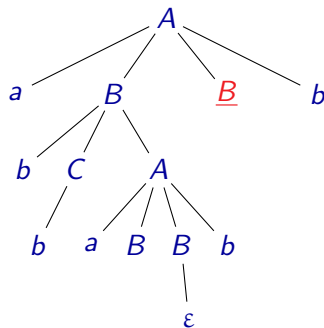
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$\underline{B} \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



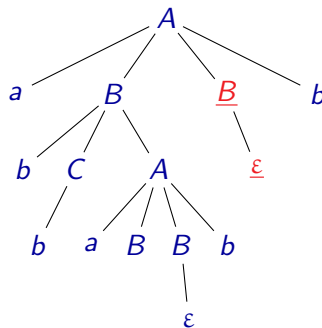
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBb\underline{B}b$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$\underline{B} \rightarrow \underline{\epsilon} \mid bCA$

$C \rightarrow AB \mid a \mid b$



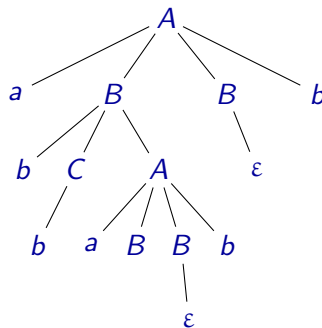
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBb\underline{B}b \Rightarrow abbaBbb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



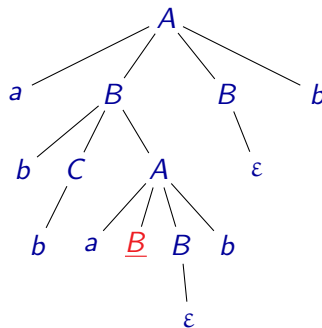
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb \Rightarrow abbaBbb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$\underline{B} \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb \Rightarrow abba\underline{B}bb$

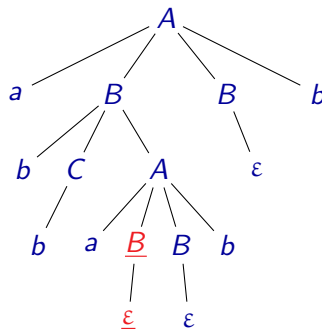


# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$\underline{B} \rightarrow \underline{\epsilon} \mid bCA$

$C \rightarrow AB \mid a \mid b$



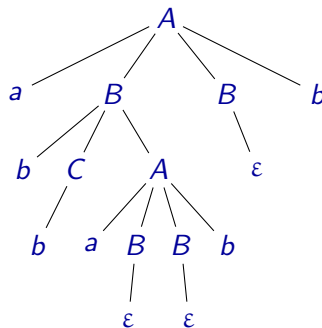
$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb \Rightarrow$   
 $abba\underline{B}bb \Rightarrow abbabb$

# Derivation Tree

$A \rightarrow aBBb \mid AaA$

$B \rightarrow \varepsilon \mid bCA$

$C \rightarrow AB \mid a \mid b$



$A \Rightarrow aBBb \Rightarrow abCABb \Rightarrow abCaBBbBb \Rightarrow abCaBbBb \Rightarrow abbaBbBb \Rightarrow abbaBbb \Rightarrow abbabb$

For each derivation there is some **derivation tree**:

- Nodes of the tree are labelled with terminals and nonterminals.
- The root of the tree is labelled with the initial nonterminal.
- The leafs of the tree are labelled with terminals or with symbols  $\epsilon$ .
- The remaining nodes of the tree are labelled with nonterminals.
- If a node is labelled with some nonterminal  $A$  then its children are labelled with the symbols from the right-hand side of some rewriting rule  $A \rightarrow \alpha$ .

# Left and Right Derivation

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

A **left derivation** is a derivation where in every step we always replace the leftmost nonterminal.

$$\underline{E} \Rightarrow \underline{E} + E \Rightarrow \underline{E} * E + E \Rightarrow a * \underline{E} + E \Rightarrow a * a + \underline{E} \Rightarrow a * a + a$$

A **right derivation** is a derivation where in every step we always replace the rightmost nonterminal.

$$\underline{E} \Rightarrow E + \underline{E} \Rightarrow \underline{E} + a \Rightarrow E * \underline{E} + a \Rightarrow \underline{E} * a + a \Rightarrow a * a + a$$

A derivation need not be left or right:

$$\underline{E} \Rightarrow \underline{E} + E \Rightarrow E * \underline{E} + E \Rightarrow E * a + \underline{E} \Rightarrow \underline{E} * a + a \Rightarrow a * a + a$$

# Left and Right Derivation

- There can be several different derivations corresponding to one derivation tree.
- For every derivation tree, there is exactly one left and exactly one right derivation corresponding to the tree.

# Equivalence of Grammars

Grammars  $\mathcal{G}_1$  and  $\mathcal{G}_2$  are **equivalent** if they generate the same language, i.e., if  $\mathcal{L}(\mathcal{G}_1) = \mathcal{L}(\mathcal{G}_2)$ .

**Remark:** The problem of equivalence of context-free grammars is algorithmically undecidable. It can be shown that it is not possible to construct an algorithm that would decide for any pair of context-free grammars if they are equivalent or not.

Even the problem to decide if a grammar generates the language  $\Sigma^*$  is algorithmically undecidable.

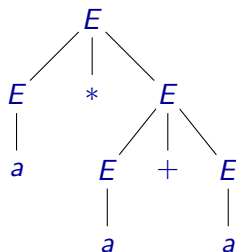
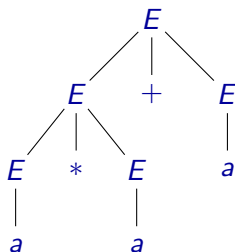
# Ambiguous Grammars

A grammar  $\mathcal{G}$  is **ambiguous** if there is a word  $w \in \mathcal{L}(\mathcal{G})$  that has two different derivation trees, resp. two different left or two different right derivations.

## Example:

$E \Rightarrow E + E \Rightarrow E * E + E \Rightarrow a * E + E \Rightarrow a * a + E \Rightarrow a * a + a$

$E \Rightarrow E * E \Rightarrow E * E + E \Rightarrow a * E + E \Rightarrow a * a + E \Rightarrow a * a + a$



# Ambiguous Grammars

Sometimes it is possible to replace an ambiguous grammar with a grammar generating the same language but which is not ambiguous.

**Example:** A grammar

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

can be replaced with the equivalent grammar

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

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

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

**Remark:** If there is no unambiguous grammar equivalent to a given ambiguous grammar, we say it is **inherently ambiguous**.



## Definition

A language  $L$  is **context-free** if there exists some context-free grammar  $\mathcal{G}$  such that  $L = \mathcal{L}(\mathcal{G})$ .

The class of context-free languages is closed with respect to:

- concatenation
- union
- iteration

The class of context-free languages is not closed with respect to:

- complement
- intersection

# Context-Free Languages

We have two grammars  $\mathcal{G}_1 = (\Pi_1, \Sigma, S_1, P_1)$  and  $\mathcal{G}_2 = (\Pi_2, \Sigma, S_2, P_2)$ , and can assume that  $\Pi_1 \cap \Pi_2 = \emptyset$  and  $S \notin \Pi_1 \cup \Pi_2$ .

- Grammar  $\mathcal{G}$  such that  $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}_1) \cdot \mathcal{L}(\mathcal{G}_2)$ :

$$\mathcal{G} = (\Pi_1 \cup \Pi_2 \cup \{S\}, \Sigma, S, P_1 \cup P_2 \cup \{S \rightarrow S_1 S_2\})$$

- Grammar  $\mathcal{G}$  such that  $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}_1) \cup \mathcal{L}(\mathcal{G}_2)$ :

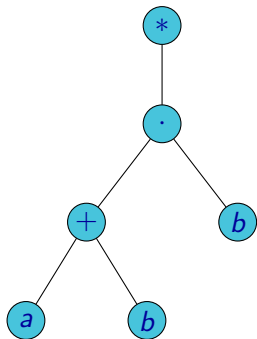
$$\mathcal{G} = (\Pi_1 \cup \Pi_2 \cup \{S\}, \Sigma, S, P_1 \cup P_2 \cup \{S \rightarrow S_1, S \rightarrow S_2\})$$

- Grammar  $\mathcal{G}$  such that  $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}_1)^*$ :

$$\mathcal{G} = (\Pi_1 \cup \{S\}, \Sigma, S, P_1 \cup \{S \rightarrow \varepsilon, S \rightarrow S_1 S\})$$

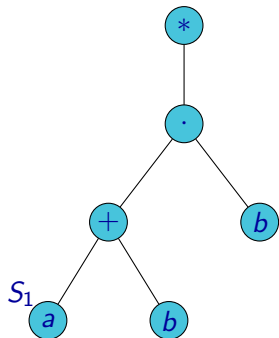
# A Context-Free Grammar for a Regular Expression

**Example:** The construction of a context-free grammar for regular expression  $((a + b) \cdot b)^*$ :



# A Context-Free Grammar for a Regular Expression

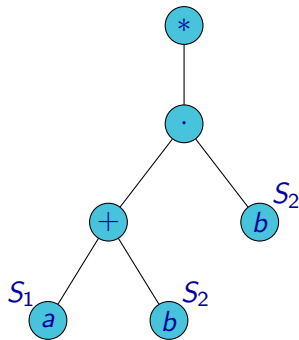
**Example:** The construction of a context-free grammar for regular expression  $((a + b) \cdot b)^*$ :



$S_1 \rightarrow a$

# A Context-Free Grammar for a Regular Expression

**Example:** The construction of a context-free grammar for regular expression  $((a + b) \cdot b)^*$ :

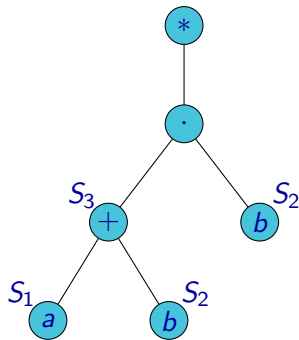


$$S_2 \rightarrow b$$

$$S_1 \rightarrow a$$

# A Context-Free Grammar for a Regular Expression

**Example:** The construction of a context-free grammar for regular expression  $((a + b) \cdot b)^*$ :



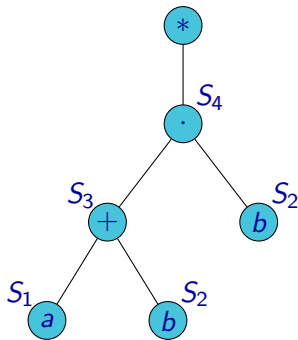
$$S_3 \rightarrow S_1 \mid S_2$$

$$S_2 \rightarrow b$$

$$S_1 \rightarrow a$$

# A Context-Free Grammar for a Regular Expression

**Example:** The construction of a context-free grammar for regular expression  $((a + b) \cdot b)^*$ :



$$S_4 \rightarrow S_3 S_2$$

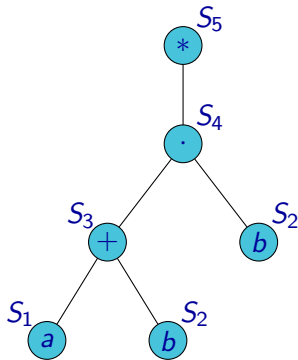
$$S_3 \rightarrow S_1 \mid S_2$$

$$S_2 \rightarrow b$$

$$S_1 \rightarrow a$$

# A Context-Free Grammar for a Regular Expression

**Example:** The construction of a context-free grammar for regular expression  $((a + b) \cdot b)^*$ :



$$S_5 \rightarrow \varepsilon \mid S_4 S_5$$

$$S_4 \rightarrow S_3 S_2$$

$$S_3 \rightarrow S_1 \mid S_2$$

$$S_2 \rightarrow b$$

$$S_1 \rightarrow a$$



# Lexical and Syntactic Analysis — an example

**Example:** We would like to recognize a language of arithmetic expressions containing expressions such as:

34

$x+1$

$-x * 2 + 128 * (y - z / 3)$

- The expressions can contain number constants — sequences of digits 0, 1, ..., 9.
- The expressions can contain names of variables — sequences consisting of letters, digits, and symbol “\_”, which do not start with a digit.
- The expressions can contain basic arithmetic operations — “+”, “-”, “\*”, “/”, and unary “-”.
- It is possible to use parentheses — “(” and “)”, and to use a standard priority of arithmetic operations.

# Lexical and Syntactic Analysis — an example

The problem we want to solve:

- **Input:** a sequence of characters (e.g., a string, a text file, etc.)
- **Output:** an abstract syntax tree representing the structure of a given expression, or an information about a syntax error in the expression

# Lexical and Syntactic Analysis — an example

It is convenient to decompose this problem into several parts:

- **Lexical analysis** — recognizing of **lexical elements** (so called **tokens**) such as for example identifiers, number constants, operators, etc.
- **Syntactic analysis** — determining whether a given sequence of tokens corresponds to an allowed structure of expressions; basically, it means finding corresponding derivation (resp. derivation tree) for a given word in a context-free grammar representing the given language (e.g., in our case, the language of all well-formed expressions).
- **Construction of an abstract syntax tree** — this phase is usually connected with the syntax analysis, where the result, actually produced by the program, is typically not directly a derivation tree but rather some kind of abstract syntax tree or performing of some actions connected with rules of the given grammar.

# Lexical and Syntactic Analysis — an example

**Terminals** for the grammar representing well-formed expressions:

- $\langle \textit{id} \textit{ent} \rangle$  — identifier, e.g. “x”, “q3”, “count\_r12”
- $\langle \textit{num} \rangle$  — number constant, e.g. “5”, “42”, “65535”
- “(” — left parenthesis
- )” — right parenthesis
- “+” — plus
- “-” — minus
- “\*” — star
- “/” — slash

**Remark:** Recognizing of sequences of symbols that correspond to individual terminals is the goal of lexical analysis.

# Lexical and Syntactic Analysis — an example

**Example:** Expression  $-x * 2 + 128 * (y - z / 3)$  is represented by the following sequence of symbols:

$- \ x \ \_ \ * \ \_ \ 2 \ \_ \ + \ \_ \ 1 \ 2 \ 8 \ \_ \ * \ \_ \ ( \ y \ \_ \ - \ \_ \ z \ \_ \ / \ \_ \ 3 \ )$

The following sequence of **tokens** corresponds to this sequence of symbols; these tokens are terminal symbols of the given context-free grammar:

$- \ \langle ident \rangle \ * \ \langle num \rangle \ + \ \langle num \rangle \ * \ ( \ \langle ident \rangle \ - \ \langle ident \rangle \ / \ \langle num \rangle \ )$

# Lexical and Syntactic Analysis — an example

The context-free grammar for the given language — the first try:

$$E \rightarrow \langle ident \rangle \mid \langle num \rangle \mid ( E ) \mid - E \mid E + E \mid E - E \mid E * E \mid E / E$$

# Lexical and Syntactic Analysis — an example

The context-free grammar for the given language — the first try:

$$E \rightarrow \langle \textit{ident} \rangle \mid \langle \textit{num} \rangle \mid ( E ) \mid - E \mid E + E \mid E - E \mid E * E \mid E / E$$

This grammar is ambiguous.

# Lexical and Syntactic Analysis — an example

The context-free grammar for the given language — the second try:

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

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

$$F \rightarrow \langle \textit{ident} \rangle \mid \langle \textit{num} \rangle \mid ( E ) \mid - F$$

Different levels of priority are represented by different nonterminals:

- $E$  — expression
- $T$  — term
- $F$  — factor

This grammar is unambiguous.



# Lexical and Syntactic Analysis — an example

The context-free grammar for the given language — the third try:

$$E \rightarrow T \mid T A E$$

$$A \rightarrow + \mid -$$

$$T \rightarrow F \mid F M T$$

$$M \rightarrow * \mid /$$

$$F \rightarrow \langle \textit{ident} \rangle \mid \langle \textit{num} \rangle \mid ( E ) \mid - F$$

We create separate nonterminals for operators on different levels of priority:

- $A$  — additive operator
- $M$  — multiplicative operator

# Lexical and Syntactic Analysis — an example

The context-free grammar for the given language — the fourth try:

$$S \rightarrow E \langle eof \rangle$$

$$E \rightarrow T \mid T A E$$

$$A \rightarrow + \mid -$$

$$T \rightarrow F \mid F M T$$

$$M \rightarrow * \mid /$$

$$F \rightarrow \langle ident \rangle \mid \langle num \rangle \mid ( E ) \mid - F$$

- It is useful to introduce special nonterminal  $\langle eof \rangle$  representing the end of input.
- Moreover, in this grammar the initial nonterminal  $S$  does not occur on the right hand side of any grammar.

# Implementation of Lexical Analysis

Enumerated type *Token\_kind* representing different kinds of **tokens**:

<b>T_EOF</b>	— the end of input
<b>T_Ident</b>	— identifier
<b>T_Number</b>	— number constant
<b>T_LParen</b>	— “(”
<b>T_RParen</b>	— “)”
<b>T_Plus</b>	— “+”
<b>T_Minus</b>	— “-”
<b>T_Star</b>	— “*”
<b>T_Slash</b>	— “/”

# Implementation of Lexical Analysis

Variable *c* : a currently processed character (resp. a special value *<eof>* representing the end of input):

- at the beginning, the first character in the input is read to variable *c*
- function `NEXT-CHAR()` returns a next character from the input

Some helper functions:

- `ERROR()` — outputs an information about a syntax error and aborts the processing of the expression
- `is-ident-start-char(c)` — tests whether *c* is a character that can occur at the beginning of an identifier
- `is-ident-normal-char(c)` — tests whether *c* is a character that can occur in an identifier (on other positions except beginning)
- `is-digit(c)` — tests whether *c* is a digit

# Implementation of Lexical Analysis

Some other helper functions:

- `CREATE-IDENT(s)` — creates an identifier from a given string *s*
- `CREATE-NUMBER(s)` — creates a number from a given string *s*

Auxiliary variables:

- *last-ident* — the last processed identifier
- *last-num* — the last processed number constant

Function `NEXT-TOKEN()` — the main part of the lexical analyser, it returns the following token from the input

# Implementation of Lexical Analysis

---

```
NEXT-TOKEN():  
  while  $c \in \{ " ", "\t" \}$  do  
     $c := \text{NEXT-CHAR}()$ ;  
  if  $c == \langle \text{eof} \rangle$  then return T_EOF  
  else switch  $c$  do  
    case "(" do  $c := \text{NEXT-CHAR}()$ ; return T_LParen  
    case ")" do  $c := \text{NEXT-CHAR}()$ ; return T_RParen  
    case "+" do  $c := \text{NEXT-CHAR}()$ ; return T_Plus  
    case "-" do  $c := \text{NEXT-CHAR}()$ ; return T_Minus  
    case "*" do  $c := \text{NEXT-CHAR}()$ ; return T_Star  
    case "/" do  $c := \text{NEXT-CHAR}()$ ; return T_Slash  
    otherwise do  
      if is-ident-start-char( $c$ ) then return SCAN-IDENT()  
      else if is-digit( $c$ ) then return SCAN-NUMBER()  
      else ERROR()
```

---

# Implementation of Lexical Analysis

---

---

SCAN-IDENT():

$s := c$

$c := \text{NEXT-CHAR}()$

**while** **is-ident-normal-char**( $c$ ) **do**

$s := s \cdot c$

$c := \text{NEXT-CHAR}()$

$\text{last-ident} := \text{CREATE-IDENT}(s)$

**return** **T\_Ident**

---

# Implementation of Lexical Analysis

---

---

```
SCAN-NUMBER ():
```

```
  s := c
```

```
  c := NEXT-CHAR()
```

```
  while is-digit(c) do
```

```
    | s := s · c
```

```
    | c := NEXT-CHAR()
```

```
  last-num := CREATE-NUMBER(s)
```

```
  return T_Number
```

---



# Implementation of Syntactic Analysis

Variable  $t$  :

- the last processed token

A helper function:

- `INIT-SCANNER()`:
  - initializes the lexical analyser
  - reads the first character from the input into variable `c`, aby tam byl nachystán pro následná volání funkce `NEXT-TOKEN()`

Reading a next token:

- `NEXT-TOKEN()`:
  - this is the previously described main function of the lexical analyser
  - by repeatedly calling this function we read the tokens
  - variable `c` always contains the symbol that has been read last

# Implementation of Syntactic Analysis

One of the often used methods of syntactic analysis is **recursive descent**:

- For each nonterminal there is a corresponding function — the function corresponding to nonterminal  $A$  implements all rules with nonterminal  $A$  on the left-hand side.
- In a given function, the next token is used to select between corresponding rules.
- Instructions in the body of a function correspond to processing of right-hand sides of the rules:
  - an occurrence of nonterminal  $B$  — the function corresponding to nonterminal  $B$  is called
  - an occurrence of terminal  $a$  — it is checked that the following token corresponds to terminal  $a$ , when it does, the next token is read, otherwise an error is reported

# Implementation of Syntactic Analysis

The previously described grammar is not very suitable for the recursive descent because it is not possible for nonterminals  $E$  and  $T$  to determine in a deterministic way one of the given pair of rules by use of just one following symbol:

$$S \rightarrow E \langle \text{eof} \rangle$$

$$E \rightarrow T \mid T A E$$

$$A \rightarrow + \mid -$$

$$T \rightarrow F \mid F M T$$

$$M \rightarrow * \mid /$$

$$F \rightarrow \langle \text{ident} \rangle \mid \langle \text{num} \rangle \mid ( E ) \mid - F$$

For example, if we want to rewrite nonterminal  $T$  and we know that the following terminal in the input is  $\langle \text{num} \rangle$ , this terminal can be generated by use of any of the rules

$$T \rightarrow F$$

$$T \rightarrow F M T$$

The following modified grammar does not have this problem:

$$S \rightarrow E \langle eof \rangle$$

$$E \rightarrow T G$$

$$G \rightarrow \varepsilon \mid A T G$$

$$A \rightarrow + \mid -$$

$$T \rightarrow F U$$

$$U \rightarrow \varepsilon \mid M F U$$

$$M \rightarrow * \mid /$$

$$F \rightarrow - F \mid ( E ) \mid \langle ident \rangle \mid \langle num \rangle$$

# Implementation of Syntactic Analysis

---

---

PARSE ():

    INIT-SCANNER()  
     $t := \text{NEXT-TOKEN}()$   
    PARSE-S()

---

$$S \rightarrow E \langle \text{eof} \rangle$$

---

---

PARSE-S ():

    PARSE-E()  
    if  $t \neq \mathbf{T\_EOF}$  then ERROR()

---

$$E \rightarrow T G$$

---

---

PARSE-E ():

    PARSE-T()  
    PARSE-G()

---

$$G \rightarrow \varepsilon \mid A T G$$

---

---

PARSE-G ():

**if**  $t \in \{\mathbf{T\_Plus}, \mathbf{T\_Minus}\}$  **then**  
        PARSE-A()  
        PARSE-T()  
        PARSE-G()

---

# Implementation of Syntactic Analysis

$$T \rightarrow F U$$

---

---

PARSE-T ():

    PARSE-F()  
    PARSE-U()

---

$$U \rightarrow \varepsilon \mid M F U$$

---

---

PARSE-U (*e1*):

**if**  $t \in \{\mathbf{T\_Star}, \mathbf{T\_Slash}\}$  **then**  
        PARSE-M()  
        PARSE-F()  
        PARSE-U()

---

# Implementation of Syntactic Analysis

$$A \rightarrow + \mid -$$

---

PARSE-A ():

```
switch  $t$  do
  case T_Plus do
    |  $t := \text{NEXT-TOKEN}()$ 
  case T_Minus do
    |  $t := \text{NEXT-TOKEN}()$ 
  otherwise do ERROR()
```

---



# Implementation of Syntactic Analysis

$$M \rightarrow * \mid /$$

---

---

PARSE-M ():

```
switch t do
| case T_Star do
|   | t := NEXT-TOKEN()
| case T_Slash do
|   | t := NEXT-TOKEN()
| otherwise do ERROR()
```

---

# Implementation of Syntactic Analysis

$$F \rightarrow \langle \text{ident} \rangle$$
$$| \langle \text{num} \rangle$$
$$| ( E )$$
$$| - F$$

PARSE-F ():

```
switch t do
  case T_Ident do
    | t := NEXT-TOKEN()
  case T_Number do
    | t := NEXT-TOKEN()
  case T_LParen do
    | t := NEXT-TOKEN()
    | PARSE-E()
    | if t ≠ T_RParen then ERROR()
    | t := NEXT-TOKEN()
  case T_Minus do
    | t := NEXT-TOKEN()
    | PARSE-F()
  otherwise do ERROR()
```

# Implementation of Syntactic Analysis

- If a function ends with a recursive call of itself, as for example function `PARSE-G()`, it is possible to replace this recursion with an iteration.
- Functions `PARSE-E()` and `PARSE-G()` can be merged into one function.
- Similarly, it is possible to replace a recursion with an iteration in function `PARSE-U()`, and functions `PARSE-T()` and `PARSE-U()` can be merged into one function.

$$E \rightarrow T G$$

$$G \rightarrow \varepsilon \mid A T G$$

---



---

PARSE-E ():

```

  PARSE-T()
  while  $t \in \{\mathbf{T\_Plus}, \mathbf{T\_Minus}\}$  do
    PARSE-A()
    PARSE-T()

```

---

$$T \rightarrow F U$$

$$U \rightarrow \varepsilon \mid M F U$$

---



---

PARSE-T ():

```

  PARSE-F()
  while  $t \in \{\mathbf{T\_Star}, \mathbf{T\_Slash}\}$  do
    PARSE-M()
    PARSE-F()

```

---

# Implementation of Syntactic Analysis

- The implementation described above just finds out whether the given input corresponds to some word that can be generated by the given grammar.
- If this is the case, it reads whole input and finishes successfully.
- If it is not the case, function `ERROR()` is called.
- In real implementation, it is useful to provide function `ERROR()` with error messages describing the kind of error together with the information about a position in the input where the error occurred (e.g., this line and column where the currently processed token starts). Function `ERROR()` can use this information to create error messages that are displayed to a user.

# Implementation of Syntactic Analysis

- Typically, we do not want to use syntactic analysis just to check that the input is correct but also to create abstract syntax tree or to perform some other types of actions connected with individual rules of the grammar.
- The previously presented code can be used as a base that can be extended with other actions such as construction of an abstract syntax tree, modifications of read expressions, and possibly some other types of computation.
- When the functions that correspond to nonterminals should create the corresponding abstract syntax tree, they can return the constructed subtree, corresponding to the part of the expression generated from the given nonterminal, as a return value.

# Implementation of Syntactic Analysis

Construction of an **abstract syntax tree**:

- An enumerated type representing binary arithmetic operations:  
`enum Bin_op { Add, Sub, Mul, Div }`
- An enumerated type representing unary arithmetic operations:  
`enum Un_op { Un_minus }`
- Functions for creation of different kinds of nodes of an abstract syntax tree:
  - `MK-VAR(ident)` — creates a leaf representing a variable
  - `MK-NUM(num)` — creates a leaf representing a number constant
  - `MK-UNARY(op, e)` — creates a node with one child *e*, on which a unary operation *op* (of type *Un\_op*) is applied
  - `MK-BINARY(op, e1, e2)` — creates a node with two children *e1* and *e2*, on which a binary operation *op* (of type *Bin\_op*) is applied

$$S \rightarrow E \langle eof \rangle$$

---

---

PARSE ():

```
  INIT-SCANNER()
  t := NEXT-TOKEN()
  e := PARSE-E()
  if t ≠ T_EOF then ERROR()
  return e
```

---



$$\begin{aligned} E &\rightarrow T G \\ G &\rightarrow \varepsilon \mid A T G \end{aligned}$$

---

---

PARSE-E ():

```
e1 := PARSE-T()  
while t ∈ {T_Plus, T_Minus} do  
    | op := PARSE-A()  
    | e2 := PARSE-T()  
    | e1 := MK-BINARY(op, e1, e2)  
return e1
```

---

# Implementation of Syntactic Analysis

$$A \rightarrow + \mid -$$

---

---

PARSE-A ():

```
switch  $t$  do
  case T_Plus do
    |  $t := \text{NEXT-TOKEN}()$ 
    | return Add
  case T_Minus do
    |  $t := \text{NEXT-TOKEN}()$ 
    | return Sub
  otherwise do ERROR()
```

---

# Implementation of Syntactic Analysis

$$\begin{aligned}T &\rightarrow F U \\ U &\rightarrow \varepsilon \mid M F U\end{aligned}$$

---

---

PARSE-T ():

```
e1 := PARSE-F()  
while t ∈ {T_Star, T_Slash} do  
    op := PARSE-M()  
    e2 := PARSE-F()  
    e1 := MK-BINARY(op, e1, e2)  
return e1
```

---

# Implementation of Syntactic Analysis

$$M \rightarrow * \mid /$$

---

---

PARSE-M():

```
switch t do
  case T_Star do
    | t := NEXT-TOKEN()
    | return Mul
  case T_Slash do
    | t := NEXT-TOKEN()
    | return Div
  otherwise do ERROR()
```

---

$$\begin{aligned}
 F &\rightarrow \langle \textit{ident} \rangle \\
 &| \langle \textit{num} \rangle \\
 &| ( E ) \\
 &| - F
 \end{aligned}$$

PARSE-F ():

```

switch t do
  case T_Ident do
    e := MK-VAR(last-ident)
    t := NEXT-TOKEN()
    return e
  case T_Number do
    e := MK-NUM(last-num)
    t := NEXT-TOKEN()
    return e
  case T_LParen do
    t := NEXT-TOKEN()
    e := PARSE-E()
    if t ≠ T_RParen then ERROR()
    t := NEXT-TOKEN()
    return e
  case T_Minus do
    t := NEXT-TOKEN()
    e := PARSE-F()
    return MK-UNARY(Un_minus, e)
  otherwise do ERROR()

```