



# Введение в синтаксический анализ

## Слова, языки и грамматики

**Автор:** Григорьев Семён

Санкт-Петербургский государственный университет  
Математико-Механический факультет  
Кафедра системного программирования

26 октября 2011г.

А. Е. Пентус, М. Р. Пентус "Теория формальных языков"

## Definition

**Алфавитом** называется конечное непустое множество. Его элементы называются *символами* (*буквами*).

## Definition

**Алфавитом** называется конечное непустое множество. Его элементы называются *символами* (*буквами*).

## Definition

**Словом** (*цепочкой, строкой*) в алфавите  $\Sigma$  называется конечная последовательность элементов  $\Sigma$ .

## Пример

- Алфавит:  $\Sigma = \{a, b, c\}$
- $a$  – слово в алфавите  $\Sigma$
- $abbc$  – слово в алфавите  $\Sigma$

## Definition

Слово, не содержащее ни одного символа (то есть последовательность длины 0), называется *пустым словом* и обозначается  $\varepsilon$ .

## Definition

Слово, не содержащее ни одного символа (то есть последовательность длины 0), называется *пустым словом* и обозначается  $\varepsilon$ .

## Definition

*Длина* слова  $\omega$ , обозначаемая  $|\omega|$ , есть число символов в  $\omega$ , причём каждый символ считается столько раз, сколько раз он встречается в  $\omega$ .

## Пример

- $|abbc| = 4$
- $|\varepsilon| = 0$



## Definition

Если  $x$  и  $y$  — слова в алфавите  $\Sigma$ , то слово  $xy$  (результат приписывания слова  $y$  в конец слова  $x$ ) называется **конкатенацией** слов  $x$  и  $y$ . Иногда конкатенацию слов  $x$  и  $y$  обозначают  $x \cdot y$ .

## Definition

Если  $x$  и  $y$  — слова в алфавите  $\Sigma$ , то слово  $xy$  (результат приписывания слова  $y$  в конец слова  $x$ ) называется **конкатенацией** слов  $x$  и  $y$ . Иногда конкатенацию слов  $x$  и  $y$  обозначают  $x \cdot y$ .

## Definition

Если  $x$  — слово и  $n \in \mathbb{N}$ , то через  $x^n$  обозначается слово  $\underbrace{x \cdot x \cdot \dots \cdot x}_{n \text{ раз}}$ .

По определению  $x^0 \Rightarrow \varepsilon$  (знак  $\Rightarrow$  читается "равно по определению").  
Всюду далее показатели над словами и символами, как правило, являются натуральными числами.

## Пример

- По принятым соглашениям:

- ▶  $ba^3 = baaa$

- ▶  $(ba)^3 = bababa$

## Definition

Говорят, что слово  $x$  – *префикс* слова  $y$  (обозначение  $x \sqsubset y$ ), если  $y = xi$  для некоторого слова  $i$ .

Пример:  $\varepsilon \sqsubset baa$ ,  $b \sqsubset baa$ ,  $ba \sqsubset baa$ ,  $baa \sqsubset baa$

## Definition

Говорят, что слово  $x$  – **префикс** слова  $y$  (обозначение  $x \sqsubset y$ ), если  $y = xi$  для некоторого слова  $i$ .

Пример:  $\varepsilon \sqsubset baa$ ,  $b \sqsubset baa$ ,  $ba \sqsubset baa$ ,  $baa \sqsubset baa$

## Definition

Говорят, что слово  $x$  – **суффикс** слова  $y$  (обозначение  $x \sqsupset y$ ), если  $y = ix$  для некоторого слова  $i$ .

## Definition

Говорят, что слово  $x$  – *подслово* (substring) слова  $y$ , если  $y = uxv$  для некоторых слов  $u$  и  $v$ .

## Definition

Через  $|w|_a$  обозначается количество вхождений символа  $a$  в слово  $w$ .

## Definition

Если  $L \subseteq \Sigma^*$ , то  $L$  называется **языком** (или **формальным языком**) над алфавитом  $\Sigma$ .

# Порождающие грамматики



# Порождающие грамматики

## Definition

**Порождающей грамматикой (грамматикой типа 0)** называется четвёрка  $G \Rightarrow \{N, \Sigma, P, S\}$ , где  $N$  и  $\Sigma$  – конечные алфавиты,  $N \cap \Sigma = \emptyset$ ,  $P \subset (N \cup \Sigma)^+ \times (N \cup \Sigma)^*$ ,  $P$  конечно и  $S \in N$ .

Здесь:

- $\Sigma$  – **основной алфавит (терминальный алфавит)**, его элементы называются **терминальными символами** или **терминалами**
- $N$  – **вспомогательный алфавит (нетерминальный алфавит)**, его элементы называются **нетерминальными символами**, **нетерминалами**
- $S$  – **начальный символ**
- Пары  $(\alpha, \beta) \in P$  называются **правилами подстановки**, просто **правилами** или **продукциями** и записываются в виде  $\alpha \rightarrow \beta$

# Порождающие грамматики. Пример

Пусть даны множества:

- $N = \{S\}$
- $\Sigma = \{a, b, c\}$
- $P = \{S \rightarrow acSbcS, cS \rightarrow \varepsilon\}$

Тогда  $(N, \Sigma, P, S)$  является порождающей грамматикой.

# Порождающие грамматики

Будем обозначать:

- элементы множества  $\Sigma$  – строчным и буквам и из начала латинского алфавита
- элементы множества  $N$  – заглавными латинскими буквами

Обычно грамматику задают в виде списка правил, подразумевая, что алфавит  $N$  составляют все заглавные буквы, встречающиеся в правилах, а алфавит  $\Sigma$  – все строчные буквы, встречающиеся в правилах. При этом правила порождающей грамматики записывают в таком порядке, что левая часть первого правила есть начальный символ  $S$ .

Для обозначения  $n$  правил с одинаковыми левыми частями  $\alpha \rightarrow \beta_1, \dots, \alpha \rightarrow \beta_n$  часто используют сокращённую запись  $\alpha \rightarrow \beta_1 | \dots | \beta_n$ .

# Порождающие грамматики

## Definition

Пусть дана грамматика  $G$ . Пишем  $\varphi \xRightarrow{G} \psi$ , если  $\varphi = \eta\alpha\theta$ ,  $\psi = \eta\beta\theta$  и  $(\alpha \rightarrow \beta) \in P$  для некоторых слов  $\alpha, \beta, \eta, \theta$  в алфавите  $N \cup \Sigma$ .

## Remark

Когда из контекста ясно, о какой грамматике идёт речь, вместо  $\xRightarrow{G}$  можно писать просто  $\Rightarrow$ .

# Порождающие грамматики. Пример

Пусть  $G = \langle S, a, b, c, S \rightarrow acSbcS, cS \rightarrow \varepsilon, S \rangle$ . Тогда  $cSacS \xRightarrow[G]{} cSa$ .

# Порождающие грамматики

## Definition

Если  $\omega_0 \xRightarrow{G} \omega_1 \xRightarrow{G} \dots \xRightarrow{G} \omega_n$ , где  $n \geq 0$ , то пишем  $\omega_0 \xRightarrow{G}^* \omega_n$  (другими словами, бинарное отношение  $\xRightarrow{G}^*$  является рефлексивным, транзитивным замыканием бинарного отношения  $\xRightarrow{G}$ , определённого на множестве  $(N \cup \Sigma)^*$ ). При этом последовательность слов  $\omega_0, \omega_1, \dots, \omega_n$  называется **выводом** (derivation) слова  $\omega_n$  из слова  $\omega_0$  в грамматике  $G$ . Число  $n$  называется **длиной** (**количеством шагов**) этого вывода.

## Remark

В частности, для всякого слова  $\omega \in (N \cup \Sigma)^*$  имеет место  $\omega \xRightarrow{G}^* \omega$  (так как возможен вывод длины 0)

# Порождающие грамматики. Пример

Пусть  $G = \langle \{S\}, \{a, b\}, \{S \rightarrow aSa, S \rightarrow b\}, S \rangle$ . Тогда  $aSa \xRightarrow[G]{*} aaaaaSaaaa$ .  
Длина этого вывода – 3.



# Порождающие грамматики

## Definition

Язык, *порождаемый грамматикой*  $G$ , – это множество  $L(G) \Leftarrow \{\omega \in \Sigma^* \mid S \xRightarrow[G]{*} \omega\}$ . Будем также говорить, что грамматика  $G$  *порождает* язык  $L(G)$ .

## Remark

Существенно, что в определение порождающей грамматики включены два алфавита –  $\Sigma$  и  $N$ . Это позволяет нам "отсеять" часть слов, получаемых из начального символа. А именно, отбрасывается каждое слово, содержащее хотя бы один символ, не принадлежащий алфавиту  $\Sigma$ .

# Порождающие грамматики. Пример

Если  $G = \langle \{S\}, \{a, b\}, \{S \rightarrow aSa, S \rightarrow bb\}, S \rangle$ , то  
 $L(G) = \{a^n bba^n \mid n \geq 0\}$ .

## Definition

Две грамматики эквивалентны, если они порождают один и тот же язык.

## Пример

Грамматика  $S \rightarrow abS, S \rightarrow a$  и грамматика  $T \rightarrow aU, U \rightarrow baU, U \rightarrow \varepsilon$  эквивалентны.

## Классы грамматик

## Definition

**Контекстной грамматикой** (контекстно-зависимой грамматикой, грамматикой непосредственно составляющих, НС-грамматикой, грамматикой типа 1) (*context-sensitive grammar, phrase-structure grammar*) называется порождающая грамматика, каждое правило которой имеет вид  $\eta A \theta \rightarrow \eta \alpha \theta$ , где  $A \in N$ ,  $\eta \in (N \cup \Sigma)^*$ ,  $\theta \in (N \cup \Sigma)^*$ ,  $\alpha \in (N \cup \Sigma)^+$ .

## Definition

**Контекстно-свободной грамматикой** (КС-грамматикой, бесконтекстной грамматикой, грамматикой типа 2) (*context-free grammar*) называется порождающая грамматика, каждое правило которой имеет вид  $A \rightarrow \alpha$ , где  $A \in N$ ,  $\alpha \in (N \cup \Sigma)^*$ .