

БУ ВО Ханты-Мансийского автономного округа – Югры  
«Сургутский государственный университет»

# Теория языков программирования и методы трансляции

## Часть 2

### ■ Формальные языки и грамматики

Гришмановский Павел Валерьевич

доцент кафедры автоматизации и компьютерных систем, к.т.н., доцент

Сургут, 2019

# Формальные языки и грамматики

**Язык** – это заданный набор символов и правил их комбинирования для записи осмысленных текстов

Набор символов языка – **алфавит**

Правила комбинирования – **синтаксис**

Смысл конструкций языка – **семантика**

**Язык программирования** – формальная **знаковая** система, предназначенная для записи компьютерных программ, которая определяет набор **лексических**, **синтаксических** и **семантических** правил, задающих внешний вид программы и действия, которые выполнит исполнитель (компьютер) под ее управлением

**Задача языка программирования** – передача команд и данных от человека к компьютеру

# Формальные языки и грамматики

---

## Синтаксис — это

- система языковых категорий, относящихся к соединениям слов и строению предложений
- набор правил, определяющих множество допустимых конструкций языка
- правило записи программ на языке программирования
- система обозначений для обмена информацией между человеком и компьютером

# Формальные языки и грамматики

## Синтаксические свойства языка

- **Самодокументируемость** (легкость чтения) – синтаксические конструкции отражают семантические особенности, структуры данных и алгоритмов становятся очевидны при просмотре текста
  - наличие избыточности и умолчаний
  - естественность форматов и идентификаторов
  - наличие свободных комментариев
- **Легкость написания** – семантика передается минимально необходимыми конструкциями
  - краткость и регулярность конструкций
  - сокращенные формы конструкций и умолчания
  - естественность форматов и идентификаторов

# Формальные языки и грамматики

## Синтаксические свойства языка

- **Верифицируемость** – возможность математически строгого доказательства корректности
  - отсутствие избыточности и умолчаний
  - регулярность конструкций
- **Легкость трансляции** – возможность построения транслятора и приемлемые затраты на трансляцию
  - отсутствие избыточности
  - регулярность конструкций
- **Однозначность** – отсутствие неопределенности толкований конструкций языка, любая конструкция имеет единственный смысл

# Формальные языки и грамматики

**Семантика** – смысловое значение (интерпретация) абстрактного синтаксиса (множества допустимых видов конструкций), представленное в терминах математически строгой формальной модели

Подходы к заданию семантики

- операционный – в терминах переходов абстрактной машины из одного состояния в другое (*категориальная абстрактная машина, SECD-машина П. Лендина, SK-машина Дэвида Тёрнера*)
- пропозиционный (деривационный) – значение конструкций языка выражается формулами, описывающими состояния объектов программы (*аксиоматический метод Хоара, метод индуктивных утверждений Флойда*)
- денотационный – смысл конструкций языка представлен в терминах абстракции функций, оперирующих состояниями программы (*теория вычислений Д. Скотта*)

# Формальные языки и грамматики

Задать язык означает

- задать алфавит — множество допустимых символов
- задать синтаксис — множество правильных программ
- задать семантику — смысл программ

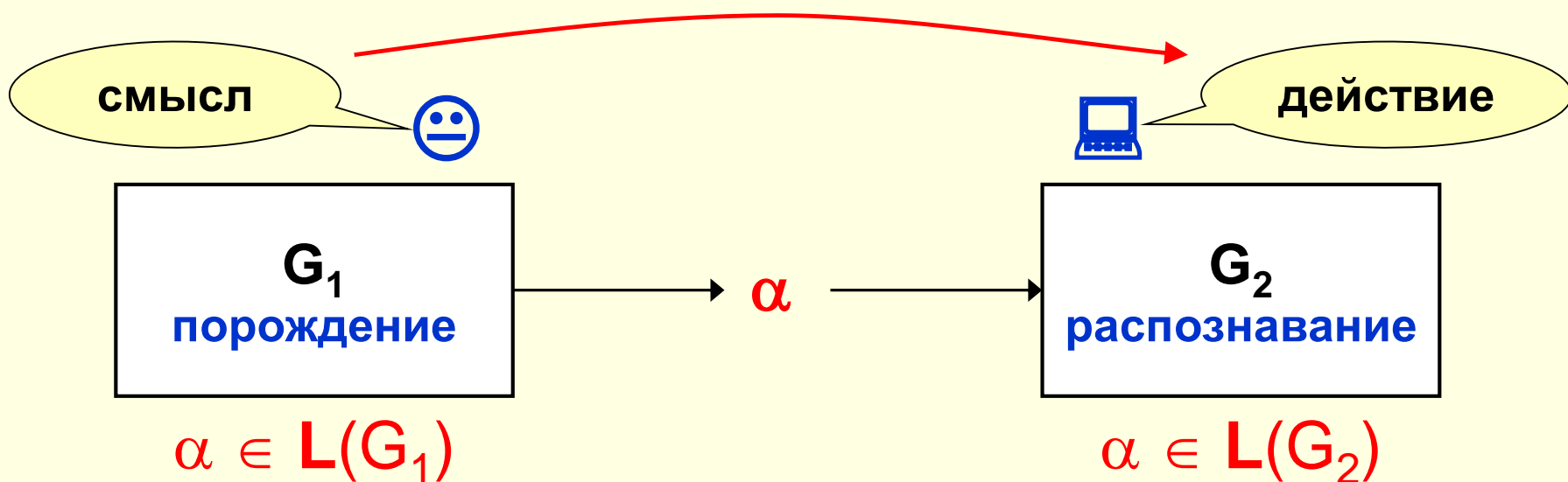
Способы задания языка

- перечисление всех допустимых цепочек языка
- определение способа порождения цепочек языка — генерация
- определение способа распознавания цепочек языка — разбор, трансляция

# Формальные языки и грамматики

**Грамматика** – математическая система, определяющая язык, т.е. формальное описание способа

- порождения цепочек символов – порождающая
- распознавания цепочек символов – распознающая





# Формальные языки и грамматики

**Язык**  $L$  как символьная система – множество допустимых цепочек  $\alpha$ , состоящих из символов алфавита  $V$ , порождаемое грамматикой  $G$

$$\alpha \in L, L = L(G(V))$$

**Алфавит**  $V$  – счетное множество допустимых символов языка, может быть бесконечным

$V^*$  – множество всех возможных цепочек над  $V$ , включая  $\lambda$

$V^+$  – множество всех возможных цепочек над  $V$ , **не** включая  $\lambda$

$$V^* = V^+ \cup \{\lambda\}, \text{ где } \lambda \text{ – пустая цепочка}$$

**Цепочка**  $\alpha(V)$  – цепочка  $\alpha$  над алфавитом  $V$ , состоящая только из символов множества  $V$

$$\forall \alpha(V) \in V^*$$

**Язык**  $L(V)$  – язык  $L$  над алфавитом  $V$ , счетное подмножество цепочек конечной длины из множества всех возможных цепочек над алфавитом  $V$ , может быть бесконечным, а длина любой цепочки может быть сколь угодно большой

$$L(V) \subseteq V^*$$

$L_1(V) \subseteq L_2(V)$ , если  $\forall \alpha \in L_1(V) : \alpha \in L_2(V)$  –  $L_2$  включает  $L_1$

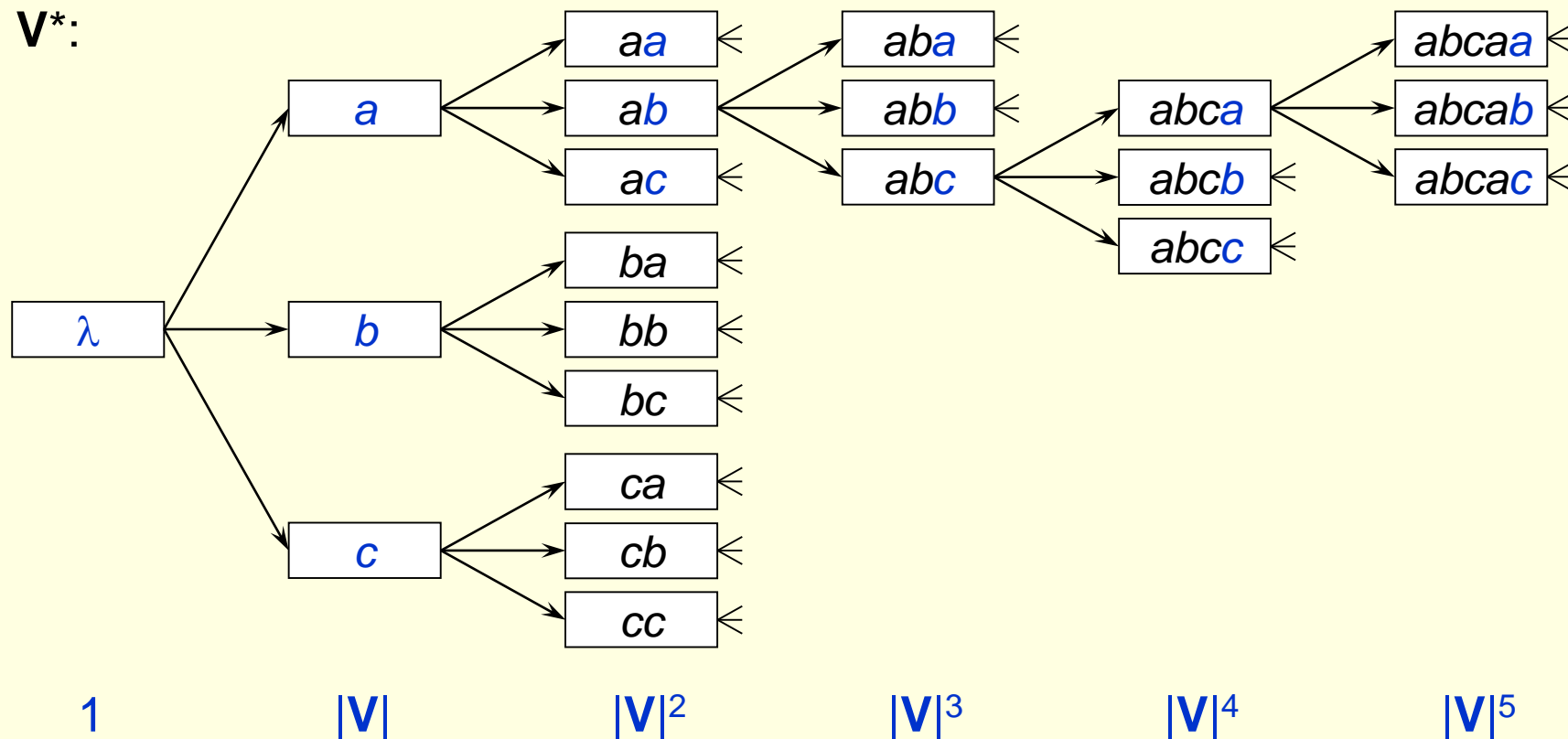
$L_1(V) = L_2(V)$ , если  $L_1(V) \subseteq L_2(V) \wedge L_2(V) \subseteq L_1(V)$  –  $L_1$  и  $L_2$  эквивалентны

$L_1(V) \cong L_2(V)$ , если  $L_1(V) \cup \{\lambda\} = L_2(V) \cup \{\lambda\}$  –  $L_1$  и  $L_2$  почти эквивалентны

# Формальные языки и грамматики

$V = \{a, b, c\}$

$V^*$ :



$V^* = \{\lambda, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, aab, aac, aba, \dots\}$

# Формальные языки и грамматики

$\alpha, \beta$	Цепочки символов	Любое количество символов, записанных друг за другом	$\alpha = a_1 a_2 a_3 \dots a_n$ $\beta = b_1 b_2 b_3 \dots b_m$
$\alpha = \beta$	Совпадение, эквивалентность, равенство	Одинаковый состав, количество и порядок символов	$a_i = b_j; i = j, n = m,$ $i \in [1, n], j \in [1, m]$
$ \alpha $	Длина цепочки	Количество символов в цепочке	$ \alpha  = n,  \beta  = m$ $ \alpha  =  \beta : \alpha = \beta$
$\alpha\beta$	Конкатенация, объединение	Дописывание символов второй цепочки к первой в том же порядке	$\alpha\beta = a_1 a_2 \dots a_n b_1 b_2 \dots b_m$ $ \alpha\beta  =  \alpha  +  \beta $ $\exists \alpha, \beta: \alpha\beta \neq \beta\alpha$
$\alpha^R$	Обращение	Запись символов цепочки строго в обратном порядке	$\alpha^R = a_n a_{n-1} a_{n-2} \dots a_2 a_1$ $ \alpha^R  =  \alpha $ $(\alpha\beta)^R = \beta^R \alpha^R$
$\alpha^k$	Итерация, повторение	Конкатенация цепочки с собой $k$ раз, $k \geq 0$	$\alpha^k = (a_1 \dots a_n)_1 \dots (a_1 \dots a_n)_k$ $ \alpha^k  = k \cdot  \alpha $

# Формальные языки и грамматики

**Пустая цепочка**  $\lambda$  ( $\varepsilon$ ,  $E$ ,  $e$ ) — это цепочка, которая не содержит ни одного символа

- $|\lambda| = 0$
- $\lambda^R = \lambda$
- $\lambda^k = \lambda, \forall k \geq 0$
- $\lambda\alpha = \alpha\lambda = \alpha, \forall \alpha \in V^*$
- $\alpha^0 = \lambda, \forall \alpha \in V^*$

# Формальные языки и грамматики

## Способы задания грамматик

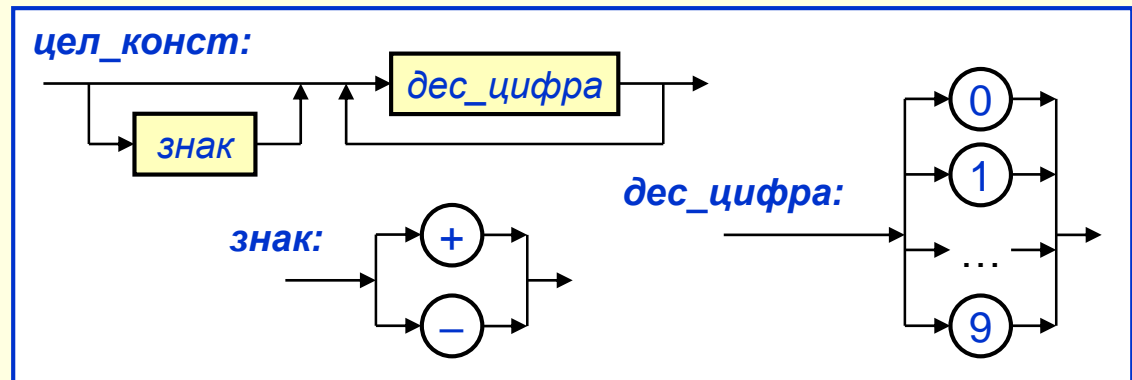
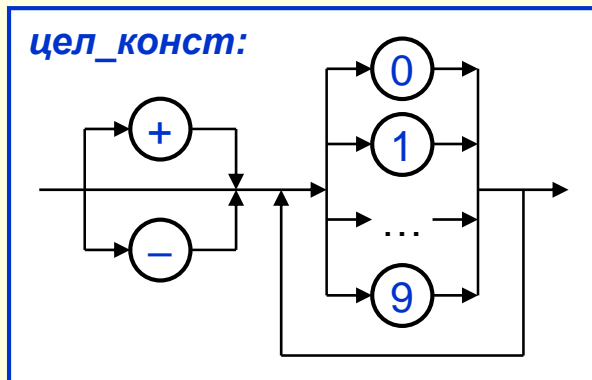
- На естественном языке – формально точное описание правил построения лексем, конструкций языка и их комбинирования
- Графически (диаграммы Вирта) – показывают порядок «сборки» конструкций языка
- Математическое определение множеств – использование формального аппарата теории множеств, алгебры логики и т.п.
- Нормальная форма Бэкуса (НФБ), форма Бэкуса-Наура (ФБН, БНФ) – формальная запись правил построения конструкций языка
- Расширенная форма Бэкуса-Наура с использованием метасимволов – аналогично, но предоставляет большую компактность, наглядность и гибкость описания
- Лямбда-исчисление – функциональная модель вычислителя на цепочках символов, эквивалент НФБ
- Конечные автоматы – простейшая формальная машина, применяемая в практике, но имеющая существенные ограничения описательной мощности
- Формальные машины – автоматы с памятью (магазинной, с произвольным доступом), в качестве управляющего устройства которых выступает конечный автомат
- Регулярные выражения – операции с цепочками символов, эквивалент конечных автоматов

# Формальные языки и грамматики

**Описание на естественном языке** – формально точное описание строения термина языка, лингвистическая формула

- *Идентификатор – это любая последовательность латинских больших и маленьких букв, символа «\_» и цифр от 0 до 9, начинающаяся не с цифры.*
- *Целочисленная константа в десятичной системе счисления представляет собой любую последовательность десятичных цифр от 0 до 9, начинающуюся не с нуля, если она не представляет значение 0, которой может предшествовать один из символов знака числа «+» или «-», и вслед за последней цифрой может быть в любой комбинации записано по одному из модификаторов «U» или «i» для обозначения беззнакового типа, и «L» или «l» для обозначения длинного целого числа.*

**Диаграммы Вирта** – ориентированный граф с вершинами двух типов: **терминальными** символами (символами алфавита языка) и **нетерминальными** символами (металингвистическими переменными, т.е. терминами языка, требующими раскрытия в цепочки терминальных символов)



# Формальные языки и грамматики

## Математическое определение множеств – формулы

- $L_1 = \{ a^n b^n : a, b \in \mathbf{V}, n > 0 \}$
- $L_2 = \{ a^n b^m c^n d^m : a, b, c, d \in \mathbf{V}, n, m \geq 0 \}$
- $L_3 = \{ a^n b^{f(n)} : a, b \in \mathbf{V}, n \geq 0, f(n) \geq 0 \text{ – некоторая (любая) вычислимая функция} \}$
- $L_B = \{ \forall \alpha : \alpha \in \{0, 1\}^*, |\alpha| = 8 \}$
- $L_C = \{ L_{\text{prep}}(\mathbf{V}) \cup L_{\text{type}}(\mathbf{V}) \cup L_{\text{var}}(\mathbf{V}) \cup L_{\text{func}}(\mathbf{V}) \}^*$

# Формальные языки и грамматики

## Регулярные выражения

1. Отдельные символы являются регулярными выражениями
  2. Если  $a$  и  $b$  – регулярные выражения, то регулярными выражениями также являются:
    - $a \vee b$  – выбор, альтернатива
    - $ab$  – конкатенация
    - $(a)$  – группировка
    - $a^*$  – повторение, замыкание Клини (0 или более раз)
  3. Ничто другое не является регулярным выражением
- $((0 \vee 1 \vee 2 \vee 3 \vee 4 \vee 5 \vee 6 \vee 7 \vee 8 \vee 9) \vee (+ \vee -)(0 \vee 1 \vee 2 \vee 3 \vee 4 \vee 5 \vee 6 \vee 7 \vee 8 \vee 9)) (0 \vee 1 \vee 2 \vee 3 \vee 4 \vee 5 \vee 6 \vee 7 \vee 8 \vee 9)^*$
  - $(a \vee b \vee c \vee d \vee e \vee f \vee g \vee h \vee i \vee j \vee k \vee l \vee m \vee n \vee o \vee p \vee q \vee r \vee s \vee t \vee u \vee v \vee w \vee x \vee y \vee z \vee \_)(a \vee b \vee c \vee d \vee e \vee f \vee g \vee h \vee i \vee j \vee k \vee l \vee m \vee n \vee o \vee p \vee q \vee r \vee s \vee t \vee u \vee v \vee w \vee x \vee y \vee z \vee \_ \vee 0 \vee 1 \vee 2 \vee 3 \vee 4 \vee 5 \vee 6 \vee 7 \vee 8 \vee 9)^*$
  - $(0 \vee 1)(0 \vee 1)(0 \vee 1)(0 \vee 1)(0 \vee 1)(0 \vee 1)(0 \vee 1)(0 \vee 1)$



# Формальные языки и грамматики

**Метасимволы** – используются самостоятельно или для сокращения записи в расширенной форме Бэкуса–Наура

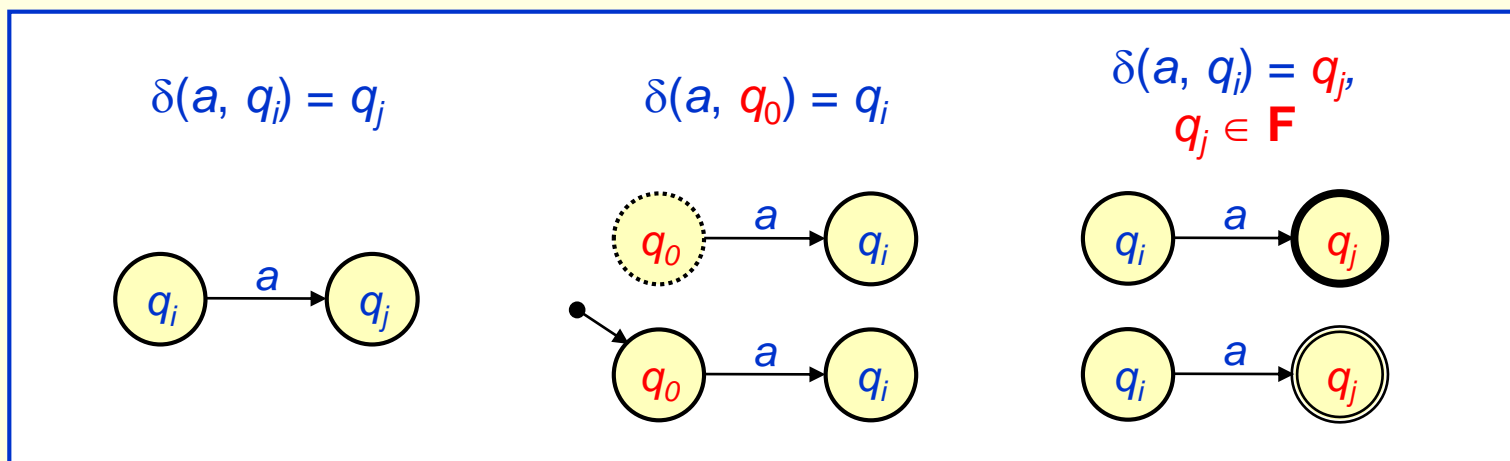
- $\langle \dots \rangle$  – нетерминальный символ (металингвистическая переменная)
- $"\dots"$  – запись терминальных символов, совпадающих с метасимволами
- $[\dots]$  – необязательная часть
- $\{\dots\}$  – повторение (0 или более раз)
- $(\dots)$  – группировка
- $\dots|\dots$  – альтернативы
  
- $[+|-](0|1|2|3|4|5|6|7|8|9)\{0|1|2|3|4|5|6|7|8|9\}$
- $\langle \text{цел\_конст} \rangle ::= [\langle \text{знак} \rangle] \langle \text{дес\_цифра} \rangle \{ \langle \text{дес\_цифра} \rangle \}$   
 $\langle \text{знак} \rangle ::= +|-$   
 $\langle \text{дес\_цифра} \rangle ::= 0|1|2|3|4|5|6|7|8|9$
- $\langle \text{идентификатор} \rangle ::= \langle \text{буква} \rangle \{ \langle \text{буква} \rangle | \langle \text{цифра} \rangle \}$   
 $\langle \text{буква} \rangle ::= a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z|_$   
 $\langle \text{цифра} \rangle ::= 0|1|2|3|4|5|6|7|8|9$
- $[\langle \text{класс\_памяти} \rangle] \langle \text{тип\_возвр\_знач} \rangle [\langle \text{согл\_вызова} \rangle]$   
 $\langle \text{имя} \rangle "(" \langle \text{список\_форм\_парам} \rangle ")" (";" | "{" \langle \text{тело\_функции} \rangle "}")$

# Формальные языки и грамматики

## Конечные автоматы

Конечный автомат  $M = (\mathbf{Q}, \mathbf{V}, \delta, q_0, \mathbf{F})$ , где

- $\mathbf{Q}$  – конечное множество состояний
- $\mathbf{V}$  – конечное множество символов (алфавит языка)
- $\delta$  – функция переходов,  $\delta : \mathbf{V} \times \mathbf{Q} \rightarrow \mathbf{Q}$   
 $\delta(a, q_i) = \mathbf{Q}_j, a \in \mathbf{V}, q_i \in \mathbf{Q}, \mathbf{Q}_j \subseteq \mathbf{Q}$ 
  - $\forall i: |\mathbf{Q}_j| \leq 1$  – **детерминированный** конечный автомат
  - $\exists i: |\mathbf{Q}_j| > 1$  – **недетерминированный** конечный автомат
- $q_0$  – начальное состояние,  $q_0 \in \mathbf{Q}$
- $\mathbf{F}$  – множество конечных состояний,  $\mathbf{F} \subseteq \mathbf{Q}$



# Формальные языки и грамматики

## Форма Бэкуса–Наура или грамматики Хомского

- Грамматика  $G = (V^T, V^N, P, S)$ , где
  - $V^T$  – множество терминальных символов, соответствующее алфавиту порождаемого языка
  - $V^N$  – множество нетерминальных символов (металингвистических переменных), используемых в записи грамматики
    - $V^T \cap V^N = \emptyset$  – множества не пересекаются
    - $V^T \cup V^N = V$  – объединение множеств образует алфавит грамматики
  - $P$  – множество правил грамматики
    - $P : \alpha \rightarrow \beta, \alpha \in V^+, \beta \in V^*$
  - $S$  – начальный (целевой) символ грамматики,  $S \in V^N$
- В левой части правила должен присутствовать хотя бы один нетерминальный символ – порождение цепочек языка не может осуществляться на основе других цепочек того же языка
- Каждый нетерминальный символ должен хотя бы раз встретиться в левой части правила – в противном случае из него нельзя вывести цепочку терминальных символов
- Правила могут (должны!) быть рекурсивными – возможен вывод бесконечного числа цепочек языка при конечном множестве правил
- Грамматика должна содержать хотя бы одно нерекурсивное правило – иначе невозможно закончить вывод цепочки языка

# Формальные языки и грамматики

$$G_1 = (V^T, V^N, P, S),$$

$$V^T = \{ 0, 1 \}$$

$$V^N = \{ A, S \}$$

$$P = \{ S \rightarrow 0A1$$

$$0A \rightarrow 00A1$$

$$A \rightarrow \lambda \}$$



01

0011

000111

00001111

...

$$G_2 = (V^T, V^N, P, S),$$

$$V^T = \{ 0, 1 \}$$

$$V^N = \{ S \}$$

$$P = \{ S \rightarrow 0S1 \mid 01 \}$$



0

12

32453

+1

+837456587493725633

-9876

-0

007

$$G_3 = (V^T, V^N, P, S),$$

$$V^T = \{ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, - \}$$

$$V^N = \{ S, T, F \}$$

$$P = \{ S \rightarrow T \mid +T \mid -T$$

$$T \rightarrow F \mid TF$$

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

