

# Лекции по теории формальных языков

## Лекция 3.

### Автоматы с магазинной памятью. Грамматики

Александр Сергеевич Герасимов

<http://gas-teach.narod.ru>

Кафедра математических и информационных технологий  
Санкт-Петербургского академического университета  
Российской академии наук.  
Весенний семестр 2010/11 учебного года

25 февраля 2011 г.

# План

1 Автоматы с магазинной памятью

2 Грамматики

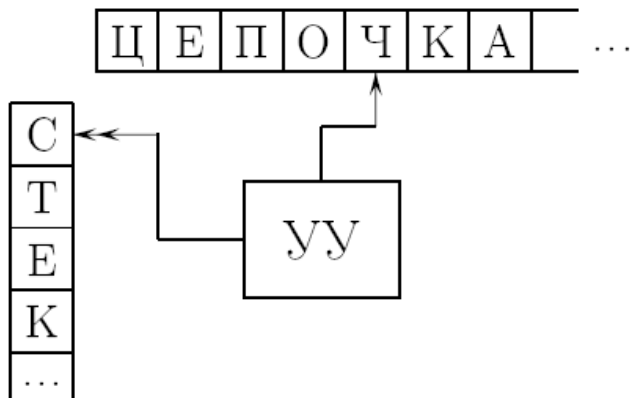
# План

1 Автоматы с магазинной памятью

2 Грамматики

# Автомат с магазинной памятью как физическое устройство

Рассмотрим более мощный, чем конечный автомат, распознаватель языка — *автомат с магазинной памятью* (МП-автомат, МПА, магазинный автомат, *pushdown automaton*, PDA).



# Неформальное описание работы МПА как физического устройства

На каждом такте устройство управления (УУ) анализирует

- своё текущее состояние  $q$ ,
- обозреваемый входной символ  $a$ ,
- символ  $B$  на вершине стека

и принимает решение о том,

- в какое состояние  $q'$  перейти,
- какую цепочку записать в стек вместо символа  $B$ ,
- следует ли сдвинуть указатель на следующую справа ячейку входной ленты ( $\rightarrow$ ) или оставить его на месте ( $\_$ ).

Описанную работу на такте мы будем рассматривать как выполнение автоматом команды

$$(q, a, B) \rightarrow (q', \gamma, \rightarrow) \quad \text{или} \quad (q, a, B) \rightarrow (q', \gamma, \_).$$

Автомат допускает цепочку, записанную на входной ленте, если он дошёл до её конца и оказался при этом в одном из заключительных состояний.

# Определение МПА

Автоматом с магазинной памятью (МПА) называется семёрка  $\mathcal{M} = (Q, \Sigma, \Gamma, \delta, q_0, F, \gamma_0)$ , где

- $Q$  — непустое конечное множество состояний,
- $\Sigma$  — входной алфавит,
- $\Gamma$  — стековый (или магазинный) алфавит,
- $\delta$  — конечное множество команд вида

$$(q, a, B) \rightarrow (q', \gamma, \rightarrow) \quad \text{или} \quad (q, a, B) \rightarrow (q', \gamma, \_)$$

(где  $q \in Q$ ,  $a \in \Sigma$ ,  $B \in \Gamma$ ,  $q' \in Q$ ,  $\gamma \in \Gamma^*$ ),

иначе говоря,  $\delta$  — конечное подмножество множества  $(Q \times \Sigma \times \Gamma) \times (Q \times \Gamma^* \times \{\rightarrow, \_ \})$ ,

- $q_0 \in Q$  — начальное состояние,
- $F \subseteq Q$  — множество заключительных состояний,
- $\gamma_0 \in \Gamma^*$  — начальное содержимое стека (магазина).

## Соглашения о символах-ограничителях

- Для удобства будем считать, что входной и стековый алфавиты МПА содержат по специальному символу-ограничителю:  $\dashv$  и  $\nabla$  соответственно.
- Символ конца входной строки  $\dashv$  размещается на входной ленте непосредственно справа от входной цепочки (и не входит в неё).
- Символ дна стека  $\nabla$  размещается под нижним символом в стеке (и в стеке больше нигде и никогда не встречается). Когда стек пуст, его верхним символом является  $\nabla$ , а при записи в стек символ дна остается на своём месте, т. е. при выполнении команды вида  $(q, a, \nabla) \rightarrow (q', \gamma, ?)$  цепочка  $\gamma$  будет записана в стек сверху символа  $\nabla$ .

# Конфигурация МПА

Конфигурацией МПА  $\mathcal{M} = (Q, \Sigma, \Gamma, \delta, q_0, F, \gamma_0)$  называется тройка  $[q, w, \gamma]$ , где

- $q \in Q$  — текущее состояние автомата,
- $w \in (\Sigma \setminus \{\#\})^*$  — суффикс входной цепочки, начинающийся с позиции указателя (и не включающий символа конца входной цепочки),
- $\gamma \in (\Gamma \setminus \{\nabla\})^*$  — содержимое стека (начиная с верхнего символа и не включая символа дна стека).



## Переход МПА в следующую конфигурацию

Пусть дан МПА  $\mathcal{M} = (Q, \Sigma, \Gamma, \delta, q_0, F, \gamma_0)$ . Для любых  $q, q' \in Q$ ,  $a \in \Sigma \setminus \{\vdash\}$ ,  $w \in (\Sigma \setminus \{\vdash\})^*$ ,  $B \in \Gamma \setminus \{\nabla\}$ ,  $\gamma, \gamma' \in (\Gamma \setminus \{\nabla\})^*$  будем писать

- $[q, aw, B\gamma] \models [q', w, \gamma'\gamma]$ , если в автомате  $\mathcal{M}$  имеется команда  $(q, a, B) \rightarrow (q', \gamma', \rightarrow)$ ;
- $[q, aw, B\gamma] \models [q', aw, \gamma'\gamma]$ , если в автомате  $\mathcal{M}$  имеется команда  $(q, a, B) \rightarrow (q', \gamma', \_)$ ;
- $[q, aw, \varepsilon] \models [q', w, \gamma']$ , если в автомате  $\mathcal{M}$  имеется команда  $(q, a, \nabla) \rightarrow (q', \gamma', \rightarrow)$ ;
- $[q, aw, \varepsilon] \models [q', aw, \gamma']$ , если в автомате  $\mathcal{M}$  имеется команда  $(q, a, \nabla) \rightarrow (q', \gamma', \_)$ ;
- $[q, \varepsilon, B\gamma] \models [q', \varepsilon, \gamma'\gamma]$ , если в автомате  $\mathcal{M}$  имеется команда  $(q, \vdash, B) \rightarrow (q', \gamma', \_)$ ;
- $[q, \varepsilon, \varepsilon] \models [q', \varepsilon, \gamma']$ , если в автомате  $\mathcal{M}$  имеется команда  $(q, \vdash, \nabla) \rightarrow (q', \gamma', \_)$ ;

и говорить, что конфигурация, указанная справа от знака  $\models$ , *следует* за конфигурацией, указанной слева от этого знака.

# Язык, распознаваемый МПА

- Пусть дан МПА  $\mathcal{M} = (Q, \Sigma, \Gamma, \delta, q_0, F, \gamma_0)$ .
- Запись  $[q, w, \gamma] \models^* [q', w', \gamma']$  будет означать, что автомат  $\mathcal{M}$  может перейти из первой конфигурации во вторую за конечное (в том числе 0) число шагов.
- МПА  $\mathcal{M}$  *распознаёт* (или *допускает*) цепочку  $w \in (\Sigma \setminus \{\vdash\})^*$ , если  $[q_0, w, \gamma_0] \models^* [f, \varepsilon, \gamma]$  для некоторых  $f \in F, \gamma \in (\Gamma \setminus \{\nabla\})^*$ .
- Множество всех цепочек, допускаемых МПА  $\mathcal{M}$ , называется *языком, распознаваемым  $\mathcal{M}$* , и обозначается  $L(\mathcal{M})$ .
- Два МПА называются *эквивалентными*, если они распознают один и тот же язык.

# Недетерминированные и детерминированные МПА

- Мы дали выше определение МПА, или недетерминированного МПА (НМПА). НМПА может иметь более одной команды с одной и той же левой частью  $(q, a, B)$ .
- Определение детерминированного МПА (ДМПА) получаем из этого определения, дополнительно потребовав, что имеется не более одной команды с одной и той же левой частью  $(q, a, B)$ , т. е. множество команд  $\delta$  — частичная функция из  $Q \times \Sigma \times \Gamma$  в  $Q \times \Gamma^* \times \{\rightarrow, \_ \}$ .
- Если не будет указано иначе, далее под МПА будет пониматься ДМПА.

# МПА, распознающий скобочный язык $LB$

- Определим МПА  $\mathcal{M} = (Q, \Sigma, \Gamma, \delta, q, F, \varepsilon)$ .
- $Q = \{q, f\}$ ,  $\Sigma = \{[, ]\}$ ,  $\Gamma = \{[\}$ ,  $F = \{f\}$ .
- Множество команд  $\delta$  будем представлять *управляющей таблицей*, проиндексированной множествами  $Q$ ,  $\Sigma$  и  $\Gamma$ : тройка индексов  $(q, a, B)$ , являющаяся левой частью команды, адресует в таблице клетку, представляющую правую часть этой команды. Не пишем в клетке состояние, если оно совпадает с текущим, и опускаем  $\_$ .

		[	]	⊣
q	[	[[, →	ε, →	
	∇	[, →		f, ε
f	[			
	∇			

$$L(\mathcal{M}) = LB.$$

## Команда допуска цепочки

- На предыдущем слайде можно было обойтись без состояния  $f$ : вместо перехода в состояние  $f$  в клетке  $(q, \uparrow, \nabla)$  достаточно было поместить указание «допустить цепочку».
- Далее мы будем использовать в МПА специальную команду допуска цепочки, записываемую как  $(q, a, B) \rightarrow \checkmark$ .
- Любой МПА с командами допуска эквивалентен некоторому МПА без них (здесь под МПА понимаются как ДМПА, так и НМПА соответственно):
  - ▶ добавим новое состояние, объявив его заключительным,
  - ▶ в этом состоянии для любых символов на входе и на стеке выполняется только сдвиг по входной цепочке,
  - ▶ заменим все команды допуска на команды перехода в новое состояние.

# МПА, распознающие язык $\{a^n b^n \mid n \geq 0\}$

1. В первом МПА состояние  $q$  — начальное и единственное заключительное. Этот МПА начинает работу с пустым стеком.

		a	b	$\vdash$
q	a	aa, $\rightarrow$	r, $\varepsilon$ , $\rightarrow$	
	$\nabla$	a, $\rightarrow$		
r	a		$\varepsilon$ , $\rightarrow$	
	$\nabla$			q, $\varepsilon$

2. Во втором МПА имеется лишь одно состояние, но введён новый стековый символ  $X$ . В начале работы в стеке находится символ  $X$ .

		a	b	$\vdash$
a			$\varepsilon$ , $\rightarrow$	
X	Xa, $\rightarrow$	$\varepsilon$	$\varepsilon$	
$\nabla$				✓

## Предложение о сравнении распознающих возможностей ДМПА и НМПА

Для произвольной цепочки  $w = a_1 \dots a_n$  через  $\overleftarrow{w}$  мы будем обозначать цепочку  $a_n \dots a_1$ .

### Предложение

*Если алфавит  $\Sigma$  содержит более одного символа, то язык  $L = \{w\overleftarrow{w} \mid w \in \Sigma^*\}$  распознаётся некоторым НМПА, но не распознаётся никаким ДМПА.*

Доказательство.

1. Определим НМПА  $M$ , распознающий язык  $L$ .

Возьмём произвольный символ  $X \notin \Sigma$ . Положим, что  $\Sigma \cup \{X\}$  — стековый алфавит,  $X$  — начальное содержимое стека.  $M$  имеет единственное состояние.

# Предложение о сравнении распознающих возможностей ДМПА и НМПА: продолжение доказательства

Управляющая таблица НМПА  $\mathcal{M}$  при  $\Sigma = \{a, b\}$ :

	a	b	$\neg$
X	$Xa, \rightarrow$ $\varepsilon$	$Xb, \rightarrow$ $\varepsilon$	$\varepsilon$
a	$\varepsilon, \rightarrow$		
b		$\varepsilon, \rightarrow$	
$\nabla$			$\checkmark$

$$L(\mathcal{M}) = L.$$



## Предложение о сравнении распознающих возможностей ДМПА и НМПА: окончание доказательства

2. Покажем (нестрого), что никакой ДМПА не распознаёт язык  $L$ .

- Предположим, что некоторый ДМПА  $M'$  распознаёт  $L$ .
- Тогда для некоторых  $w \in \Sigma^*$  и  $a \in \Sigma$  автомат  $M'$ , анализируя входную цепочку  $wa\overleftarrow{w}$ , после прочтения её префикса  $wa$ , начинает сравнивать прочитанную и непрочитанную части цепочки  $wa\overleftarrow{w}$ .
- Пусть  $M'$  после прочтения префикса  $wa$  цепочки  $wa\overleftarrow{w}$  находится в состоянии  $q$  с символом  $B$  на вершине стека.
- ДМПА  $M'$  имеет единственную команду с левой частью  $(q, a, B)$ . Эта команда указывает на начало сравнения прочитанной и непрочитанной частей цепочки  $wa\overleftarrow{w}$ .
- Но  $M'$  не сможет распознать цепочку  $wa\overleftarrow{w} \in L$ , так как начнёт сравнение сразу после прочтения префикса  $wa$  этой цепочки. Противоречие.

# МПА, распознающий язык списков $LL$

- Алфавит языка  $LL$ :  $\{a, ;, [, ]\}$ .  
Определение языка  $LL$ :
  - ▶  $a \in LL$ ;
  - ▶ если  $u \in LL$  и  $v \in LL$ , то  $[u] \in LL$  и  $u;v \in LL$ .
- МПА имеет единственное состояние, в начале работы в стеке находится символ  $L$ .

	a	;	[	]	⊖
L	$M, \rightarrow$		$L]M, \rightarrow$		
M		$L, \rightarrow$		$\epsilon$	$\epsilon$
]				$\epsilon, \rightarrow$	
∇					✓

# МПА, распознающий язык списков $LL$ : протокол обработки цепочки

Такт	Позиция указателя	Содержимое стека
1	$\diamond a; [a; a; a] \vdash$	$L\bar{\nabla}$
2	$a\diamond; [a; a; a] \vdash$	$M\bar{\nabla}$
3	$a; \diamond[a; a; a] \vdash$	$L\bar{\nabla}$
4	$a; [\diamond a; a; a] \vdash$	$L]M\bar{\nabla}$
5	$a; [a\diamond; a; a] \vdash$	$M]M\bar{\nabla}$
6	$a; [a; \diamond a; a] \vdash$	$L]M\bar{\nabla}$
7	$a; [a; a\diamond; a] \vdash$	$M]M\bar{\nabla}$
8	$a; [a; a; \diamond a] \vdash$	$L]M\bar{\nabla}$
9	$a; [a; a; a\diamond] \vdash$	$M]M\bar{\nabla}$
10	$a; [a; a; a\diamond] \vdash$	$]M\bar{\nabla}$
11	$a; [a; a; a]\diamond \vdash$	$M\bar{\nabla}$
12	$a; [a; a; a]\diamond \vdash$	$\bar{\nabla}$

# МПА, распознающий язык арифметических выражений $LA$

- Алфавит языка  $LA$ :  $\{x, +, *, (, )\}$ .  
 Определение языка  $LA$ :
  - ▶  $x \in LA$ ;
  - ▶ если  $u \in LA$  и  $v \in LA$ , то  $(u) \in LA$ ,  $u + v \in LA$  и  $u * v \in LA$ .
- МПА имеет единственное состояние, в начале работы в стеке находится символ  $E$ .

	$x$	$+$	$*$	$($	$)$	$\vdash$
$E$	$TE'$			$TE'$		
$E'$		$TE', \rightarrow$			$\epsilon$	$\epsilon$
$T$	$FT'$			$FT'$		
$T'$		$\epsilon$	$FT', \rightarrow$		$\epsilon$	$\epsilon$
$F$	$\epsilon, \rightarrow$			$E), \rightarrow$		
$)$					$\epsilon, \rightarrow$	
$\nabla$						✓

1	$\diamond(x + x) * x \vdash$	$E\Delta$
2	$\diamond(x + x) * x \vdash$	$TE'\Delta$
3	$\diamond(x + x) * x \vdash$	$FT'E'\Delta$
4	$(\diamond x + x) * x \vdash$	$E)T'E'\Delta$
5	$(\diamond x + x) * x \vdash$	$TE')T'E'\Delta$
6	$(\diamond x + x) * x \vdash$	$FT'E')T'E'\Delta$
7	$(x \diamond + x) * x \vdash$	$T'E')T'E'\Delta$
8	$(x \diamond + x) * x \vdash$	$E')T'E'\Delta$
9	$(x + \diamond x) * x \vdash$	$TE')T'E'\Delta$
10	$(x + \diamond x) * x \vdash$	$FT'E')T'E'\Delta$
11	$(x + x \diamond) * x \vdash$	$T'E')T'E'\Delta$
12	$(x + x \diamond) * x \vdash$	$E')T'E'\Delta$
13	$(x + x \diamond) * x \vdash$	$)T'E'\Delta$
14	$(x + x) \diamond * x \vdash$	$T'E'\Delta$
15	$(x + x) * \diamond x \vdash$	$FT'E'\Delta$
16	$(x + x) * x \diamond \vdash$	$T'E'\Delta$
17	$(x + x) * x \diamond \vdash$	$E'\Delta$
18	$(x + x) * x \diamond \vdash$	$\Delta$

# План

1 Автоматы с магазинной памятью

2 Грамматики

## Определение грамматики

Грамматики (и особенно контекстно-свободные грамматики) – основной способ задания формальных языков.

*Грамматикой* называется четвёрка  $G = (\Sigma, \Gamma, P, S)$ , где

- $\Sigma$  — *основной* (или *терминальный*) алфавит, элементы которого называют *терминалами* (или *терминальными символами*);
- $\Gamma$  — *вспомогательный* (или *нетерминальный*) алфавит, элементы которого называют *нетерминалами* (*нетерминальными* или *вспомогательными символами*);  $\Sigma \cap \Gamma = \emptyset$ ;
- $S \in \Gamma$  — выделенный нетерминал, называемый *аксиомой* (или *начальным нетерминалом*);
- $P$  — конечное множество слов вида  $\alpha \rightarrow \beta$  (где  $\alpha \in (\Sigma \cup \Gamma)^* \Gamma (\Sigma \cup \Gamma)^*$ ,  $\beta \in (\Sigma \cup \Gamma)^*$ ), каждое такое слово называют *правилом вывода* (просто *правилом* или *продукцией*);  $\rightarrow \notin \Sigma \cup \Gamma$ .

Пример. Грамматика  $G_1 = (\Sigma, \Gamma, P, S)$ :  $\Sigma = \{a, b\}$ ,  $\Gamma = \{S, D\}$ ,  $S$  — аксиома,  $P = \{S \rightarrow aSDa, S \rightarrow aba, aD \rightarrow Da, bD \rightarrow bb\}$ .

# Непосредственная выводимость в грамматике

Пусть  $G = (\Sigma, \Gamma, P, S)$  — грамматика,  $\gamma, \delta \in (\Sigma \cup \Gamma)^*$ .

- Говорят, что цепочка  $\delta$  непосредственно выводима в  $G$  из цепочки  $\gamma$  (и обозначают это через  $\gamma \Rightarrow_G \delta$  или через  $\gamma \Rightarrow \delta$ ), если существуют такие цепочки  $\alpha, \beta, \gamma_1, \gamma_2 \in (\Sigma \cup \Gamma)^*$ , что  $\gamma = \gamma_1 \alpha \gamma_2$ ,  $\delta = \gamma_1 \beta \gamma_2$  и  $\alpha \rightarrow \beta \in P$ .
- Пример непосредственной выводимости в  $G_1$ :
  - ▶  $\underbrace{a}_{\gamma_1} S \underbrace{Da}_{\gamma_2} \Rightarrow \underbrace{a}_{\gamma_1} aba \underbrace{Da}_{\gamma_2}$  (правило  $S \rightarrow aba \in P$ ),
  - ▶  $aDbSa \Rightarrow DabSa$  ( $aD \rightarrow Da \in P$ ),
  - ▶  $S \Rightarrow aba$  ( $S \rightarrow aba \in P$ ).



## Выводимость в грамматике

- Говорят, что цепочка  $\delta$  выводима в  $G$  из цепочки  $\gamma$  (и обозначают это через  $\gamma \Rightarrow_G^* \delta$  или через  $\gamma \Rightarrow^* \delta$ ), если существуют такие число  $n \geq 0$  и цепочки  $\eta_0, \eta_1, \dots, \eta_n$ , что

$$\gamma = \eta_0 \Rightarrow \eta_1 \Rightarrow \eta_2 \Rightarrow \dots \Rightarrow \eta_n = \delta.$$

Последовательность цепочек  $\eta_0, \eta_1, \dots, \eta_n$  называется *выводом* (в  $G$  цепочки  $\delta$  из цепочки  $\gamma$ ), а число  $n$  — *длиной* этого вывода.

- Отношение выводимости  $\Rightarrow_G^*$  есть рефлексивно-транзитивное замыкание отношения непосредственной выводимости  $\Rightarrow_G$ .

Транзитивное замыкание отношения непосредственной выводимости  $\Rightarrow_G$  мы будем обозначать через  $\Rightarrow_G^+$  (или через  $\Rightarrow^+$ ).

- Пример. Рассматриваем грамматику  $G_1$  с правилами  $P = \{S \rightarrow aSDa, S \rightarrow aba, aD \rightarrow Da, bD \rightarrow bb\}$ .  
 $S \Rightarrow_{G_1}^* a^2b^2a^2$ , поскольку

$$S \Rightarrow aSDa \Rightarrow aabaDa \Rightarrow aabDaa \Rightarrow aabbaa .$$

## Язык, порождаемый грамматикой

- Языком, порождаемым грамматикой  $G = (\Sigma, \Gamma, P, S)$ , называется множество  $L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$ .
- Грамматики  $G$  и  $G'$  называются эквивалентными, если  $L(G) = L(G')$ .
- Пример.  $aba \in L(G_1)$ ,  $a^2b^2a^2 \in L(G_1)$ . Можно показать, что  $L(G_1) = \{a^n b^n a^n \mid n \geq 1\}$ .  
Установим, что  $L(G_1) \supseteq \{a^n b^n a^n \mid n \geq 1\}$ :

- ▶ применим  $(n - 1)$  раз правило  $S \rightarrow aSDa$ , а затем правило  $S \rightarrow aba$ :

$$S \Rightarrow aSDa \Rightarrow aaSDaDa \Rightarrow \dots \Rightarrow a^{n-1}S(Da)^{n-1} \Rightarrow a^nba(Da)^{n-1};$$

- ▶ далее, пока это возможно, применяем правило  $aD \rightarrow Da$ :

$$a^nba(Da)^{n-1} \Rightarrow^* a^n b D^{n-1} a^n;$$

- ▶ наконец,  $(n - 1)$  раз применяем правило  $bD \rightarrow bb$ :

$$a^n b D^{n-1} a^n \Rightarrow^* a^n b^n a^n.$$

# Соглашение об альтернативах

- Правила  $\alpha \rightarrow \beta_1, \alpha \rightarrow \beta_2, \dots, \alpha \rightarrow \beta_n$  (называемые альтернативами) будем записывать в виде  $\alpha \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$ .
- Пример. Грамматика  $G_2 = (\Sigma, \Gamma, P, S)$ , порождающая все непустые цепочки из строчных латинских букв и десятичных цифр, начинающиеся с буквы:
  - ▶  $\Sigma = \{a, b, \dots, z, 0, 1, \dots, 9\}$ ,
  - ▶  $\Gamma = \{\langle \text{имя} \rangle, \langle \text{буква} \rangle, \langle \text{цифра} \rangle\}$ ,
  - ▶  $\langle \text{имя} \rangle$  — аксиома,
  - ▶ множество правил  $P$ :  
 $\langle \text{имя} \rangle = \langle \text{буква} \rangle | \langle \text{имя} \rangle \langle \text{буква} \rangle | \langle \text{имя} \rangle \langle \text{цифра} \rangle$ ,  
 $\langle \text{буква} \rangle = a | b | \dots | z$ ,  
 $\langle \text{цифра} \rangle = 0 | 1 | \dots | 9$ .

## Соглашения о символах и цепочках

Граматику можно задавать только множеством правил вывода, если

- терминалы обозначать строчными буквами из начала латинского алфавита ( $a, b, c, d, e, f, \dots$ );
- нетерминалы — заглавными буквами из начала латинского алфавита ( $A, B, C, D, E, F, \dots, S, T$ ) или словами в угловых скобках (например,  $\langle \text{имя} \rangle, \langle \text{буква} \rangle, \langle \text{цифра} \rangle$ );
- символы из объединения терминального и нетерминального алфавитов — заглавными буквами из конца латинского алфавита ( $U, V, W, X, Y, Z$ );
- цепочки терминалов — строчными буквами из конца латинского алфавита ( $u, v, w, x, y, z$ );
- цепочки символов из объединения терминального и нетерминального алфавитов — строчными греческими буквами ( $\alpha, \beta, \gamma, \dots, \omega$ );
- в левой части первого правила вывода стоит аксиома.

При этом можно добавлять индексы.

# Грамматика, порождающая язык $\{a^n b^n c^n \mid n \geq 1\}$

- Применим  $(n - 1)$  раз правило  $S \rightarrow aSBC$ , а затем правило  $S \rightarrow abC$ :

$$S \Rightarrow aSBC \Rightarrow aaSBCBC \Rightarrow \dots \Rightarrow a^{n-1}S(BC)^{n-1} \Rightarrow a^n b C (BC)^{n-1}.$$

- Далее, пока это возможно, применяем правило  $CB \rightarrow BC$ :

$$a^n b C (BC)^{n-1} \Rightarrow^* a^n b B^{n-1} C^n.$$

- Теперь  $(n - 1)$  раз применяем правило  $bB \rightarrow bb$ :

$$a^n b B^{n-1} C^n \Rightarrow^* a^n b^n C^n.$$

- Наконец, применяем правило  $bC \rightarrow bc$ , а затем  $(n - 1)$  раз правило  $cC \rightarrow cc$ :

$$a^n b^n C^n \Rightarrow a^n b^n c C^{n-1} \Rightarrow^* a^n b^n c^n.$$

Грамматика, порождающая язык

$$L = \{w \in \{a, b\}^* \mid \text{равны числа вхождений } a \text{ и } b \text{ в } w \}$$

- Правила  $S \rightarrow ASB \mid \varepsilon$  позволяют получить цепочки вида  $A^n B^n$  ( $n \geq 0$ ).
- С помощью правила  $AB \rightarrow BA$  можно переставлять буквы  $A$  и  $B$ .
- Наконец, правила  $A \rightarrow a$ ,  $B \rightarrow b$  превращают цепочки нетерминалов в цепочки терминалов языка  $L$ .

## Вторая грамматика, порождающая язык

$$L = \{w \in \{a, b\}^* \mid \text{равны числа вхождений } a \text{ и } b \text{ в } w\}$$

- Обозначим через  $\#_c(z)$  число вхождений символа  $c$  в цепочку  $z$ .
- Если  $au \in L$ , то  $\#_a(u) + 1 = \#_b(u)$ .  
Пусть  $B$  — нетерминал, из которого выводятся все цепочки  $u$  такие, что  $\#_a(u) + 1 = \#_b(u)$ .  
Имеем правило  $S \rightarrow aB$ .
- Аналогично, пусть  $A$  — нетерминал, из которого выводятся все цепочки  $v$  такие, что  $\#_a(v) = \#_b(v) + 1$ . Имеем правило  $S \rightarrow bA$ .
- Рассмотрим цепочку  $u$ , выводимую из  $B$ .
  - ▶ Если  $u = ax$ , то  $\#_a(x) + 2 = \#_b(x)$ . Добавим правило  $B \rightarrow aBV$ .
  - ▶ Если  $u = bx$ , то  $\#_a(x) = \#_b(x)$ . Добавим правило  $B \rightarrow bS$ .
- Аналогично будем иметь правила  $A \rightarrow bAA \mid aS$ .
- Наконец, добавим правило  $S \rightarrow \varepsilon$ .

## Грамматика, порождающая язык

$$\{wsw \mid w \in \{a, b\}^*\}$$

- С помощью правил  $S \rightarrow aAS$ ,  $S \rightarrow bBS$  и  $S \rightarrow c$  породим цепочку вида  $Vc$ , где в  $V$  за каждым  $a$  стоит  $A$  и за каждым  $b$  стоит  $B$ .
- Затем переместим  $A$  и  $B$  направо (сохранив их порядок относительно друг друга) правилами  $Aa \rightarrow aA$ ,  $Ab \rightarrow bA$ ,  $Ba \rightarrow aB$ ,  $Bb \rightarrow bB$ . Получим цепочку вида  $wWc$ , где  $W$  есть цепочка  $w$ , но записанная заглавными буквами.
- Наконец, превратим  $A$  и  $B$  в  $a$  и  $b$  соответственно и поставим их после символа  $c$  посредством правил  $Ac \rightarrow ca$ ,  $Bc \rightarrow cb$ .



# Классы грамматик и языков: иерархия Хомского

Грамматика  $G = (\Sigma, \Gamma, P, S)$  называется

- (0) грамматикой *общего вида* (грамматикой типа 0 или грамматикой без ограничений);
- (1) *контекстно-зависимой* (грамматикой типа 1 или *неукорачивающей* грамматикой), если каждое её правило имеет вид
  - ▶  $\alpha A \beta \rightarrow \alpha \gamma \beta$ , где  $\gamma \neq \varepsilon$ , или
  - ▶  $S \rightarrow \varepsilon$ , но тогда  $S$  не входит в правую часть никакого правила;
- (2) *контекстно-свободной* (грамматикой типа 2 или *бесконтекстной* грамматикой), если каждое её правило имеет вид  $A \rightarrow \alpha$ ;
- (3) *праволинейной* (грамматикой типа 3, *автоматной* или *регулярной* грамматикой), если каждое её правило имеет вид  $A \rightarrow aB$ ,  $A \rightarrow a$  или  $A \rightarrow \varepsilon$ .

Языку приписывается тип грамматики, которой он порождается.

Например, контекстно-свободные грамматики порождают так называемые контекстно-свободные языки (языки типа 2).

## Соотношения между классами языков

Пусть  $\mathbb{L}_i$  — класс языков типа  $i$  ( $i = 0, 1, 2, 3$ ).

$\mathbb{L}_3 \subsetneq \mathbb{L}_2 \subsetneq \mathbb{L}_1 \subsetneq$  класс разрешимых языков  $\subsetneq$

$\mathbb{L}_0 =$  класс перечислимых языков.

- Язык  $L$  в алфавите  $\Sigma$  называется *разрешимым*, если существует алгоритм (например, машина Тьюринга), который для любой цепочки  $w \in \Sigma^*$  выдаёт 1, если  $w \in L$ , и выдаёт 0, если  $w \notin L$ .
- Язык  $L$  в алфавите  $\Sigma$  называется *перечислимым*, если существует алгоритм, который для любой цепочки  $w \in \Sigma^*$  выдаёт 1, если  $w \in L$ , и не заканчивает работу в противном случае.

(Подробнее о разрешимых и перечислимых множествах можно прочесть, например, в книге: Герасимов А. С. Курс математической логики и теории вычислимости: Учебное пособие. СПб.: Изд-во «ЛЕМА», 2011; электронный вариант книги — на <http://gas-teach.narod.ru/cm1ct>.)

# Связь классов языков и распознавателей

- $\mathbb{L}_2$  совпадает с классом языков, распознаваемых НМПА (увидим в дальнейшем).
- $\mathbb{L}_3$  совпадает с классом языков, распознаваемых конечными автоматами (упражнение).

Практически всегда для описания синтаксиса языков программирования используются контекстно-свободные грамматики, которые мы и будем изучать.

# Литература

## Основная литература

- Замятин А. П., Шур А. М. Языки, грамматики, распознаватели: Учебное пособие. Екатеринбург : Изд-во Урал. ун-та, 2007 (электронный вариант книги — на <http://elar.usu.ru>, поиск).

## Дополнительная литература

- Ахо А., Лам М., Сети Р., Ульман Дж. Компиляторы: принципы, технологии и инструментарий. М.: ООО "И.Д. Вильямс", 2008.
- Ахо А., Ульман Дж. Теория синтаксического анализа, перевода и компиляции. М.: Мир, 1978.
- Мартыненко Б. К. Языки и трансляции: Учеб. пособие. СПб.: Издательство С.-Петербургского университета, 2004 (электронный вариант книги — на <http://www.math.spbu.ru/user/mbk>).