

# ГЛАВА 3

## Лексический анализ

- 3.1 Назначение и необходимость фазы лексического анализа
- 3.2 Транслитератор
- 3.3 Регулярные языки и грамматики
- 3.4 Конечные автоматы
  - 3.4.1 Определение КА.
  - 3.4.2 Диаграмма переходов (состояний)
- 3.5 Матрица переходов КА
- 3.6 Детерминированный и недетерминированный автомат
- 3.7 Лексический анализатор
- 3.8 Связь регулярных грамматик КА
  - 3.8.1 Построение КА на основе леволинейной грамматики
  - 3.8.2 Построение леволинейной грамматики на основе КА

## 3.1 Назначение и необходимость фазы лексического анализа

- Лексический анализ – первая фаза процесса трансляции, предназначенная для группировки символов входной цепочки в более крупные конструкции, называемые лексемами.
- Лексемы – минимальные несущие смысл объединения символов. С каждой лексемой связано два понятия:
  - класс лексемы, определяющий общее название для категории элементов, обладающих общими свойствами
  - значение лексемы, определяющее подстроку символов входной цепочки, соответствующую распознанному классу лексемы. В зависимости от класса, значение лексемы может быть преобразовано во внутреннее представление уже на этапе лексического анализа.

## *3.1 Назначение и необходимость фазы лексического анализа*

Задачи, решаемые сканером (преимущества сканера):

- представление элементарных конструкций языка в более удобной внутренней форме
- уменьшение длины программы, за счет устранения из нее несущественных для дальнейшего анализа пробелов, комментариев, игнорируемых символов.
- один и тот же язык программирования может иметь различные внешние представления элементарных конструкций.

## *3.2 Транслитератор*

- ❑ Устройство, осуществляющее сопоставление класс с каждым отдельным символом, называется **транслитератором**.
- ❑ Наиболее типичными классами символов являются:
  - буква;
  - цифра;
  - разделитель;
  - игнорируемый;
  - запрещенный;
  - прочие.

## 3.3 Регулярные языки и грамматики

*Пример. Классы лексем:*

- идентификаторы;
- служебные слова (множество идентификаторов);
- целые числа, вещественные, строки (литералы);
- однолитерные разделители ('+', '-', '(', ')'...);
- двухлитерные разделители ('//', '/\*', '\*/', ':=');
- комментарии.

## 3.3 Регулярные языки и грамматики

$\langle \text{идентификатор} \rangle ::= \text{буква} \mid \langle \text{идентификатор} \rangle \text{ буква} \mid$   
 $\langle \text{идентификатор} \rangle \text{ цифра}$   
 $\langle \text{целое} \rangle ::= \text{цифра} \mid \langle \text{целое} \rangle \text{ цифра}$   
 $\langle \text{разделитель} \rangle ::= + \mid - \mid / \mid ( \mid ) \dots$   
 $\langle \text{разделитель} \rangle ::= \langle \text{slash} \rangle / \mid \langle \text{slash} \rangle * \mid \langle \text{ast} \rangle / \mid \langle \text{colon} \rangle =$   
 $\langle \text{ast} \rangle ::= *$   
 $\langle \text{colon} \rangle ::= :$

### Регулярные грамматики

$A \rightarrow B\gamma \mid \gamma$ , где  $A, B \in VN$ ,  $\gamma \in VT^*$

$A \rightarrow \gamma B \mid \gamma$ , где  $A, B \in VN$ ,  $\gamma \in VT^*$

### Автоматные грамматики

$A \rightarrow Bt \mid t$ , где  $A, B \in VN$ ,  $t \in VT$

$A \rightarrow tB \mid t$ , где  $A, B \in VN$ ,  $t \in VT$

## 3.3 Регулярные языки и грамматики

- Доказано, что класс регулярной и автоматной грамматики почти эквиваленты. Любую регулярную грамматику можно привести к автоматному виду.
- **Домашнее задание:**  
**Алгоритм преобразования регулярной грамматики к автоматному виду**

## ГЛАВА 3. Лексический анализ

### **3.4 КОНЕЧНЫЕ АВТОМАТЫ**



## 3.4.1 Определение КА

**Конечный автомат (КА)** – это пятерка  $(Q, V, \delta, q_0, F)$ , где:

- $Q$  – конечное множество состояний;
- $V$  – конечное множество допустимых входных символов (алфавит);
- $\delta$  – отображение множества  $Q \times VT \rightarrow Q$ , определяющее поведение автомата; отображение  $\delta$  часто называют функцией переходов;
- $q_0 \in Q$  – начальное состояние;
- $F \in Q$  – непустое множество конечных состояний.

## 3.4.1 Определение КА

- Конечный автомат называется **полностью определенным**, если в каждом его состоянии определена функция перехода для каждого входного символа.

Для  $a \in V$  и  $q \in Q$  определена  $\delta(a, q) = R, R \subseteq Q$

- Множество цепочек, допускаемых конечным автоматом, составляет определяемый им язык.
- Два конечных автомата называются **эквивалентными**, если они определяют один и тот же язык.

## 3.4.2 Диаграмма переходов (состояний)

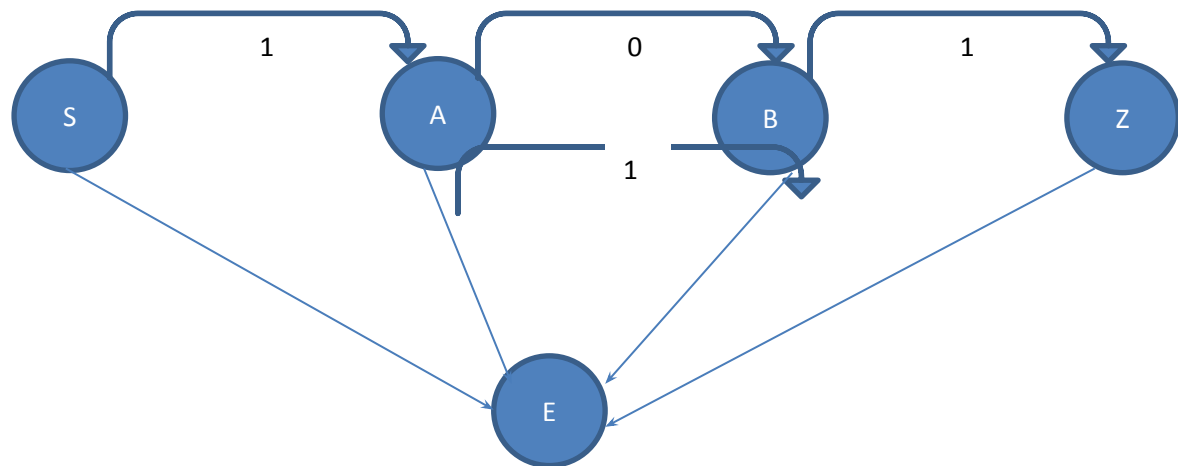
- Граф переходов (дерево, диаграмма) – направленный помеченный граф с символами состояний конечного автомата в вершинах, в котором есть дуга  $(p,q)$ , помеченная символом  $a \in V$  ( $p,q \in Q$ ), если в КА определена функция  $\delta(a,p)$  и  $q \in \delta(a,p)$ .

$M(\{S,A,B,Z\}, \{0,1\}, \delta, S, \{Z\})$

$\delta(1,S) = \{A\}$

$\delta(0,A) = \{B\}$

$\delta(1,B) = \{A,Z\}$



## **3.5 Матрица переходов КА**

- ❑ Каждая строка этой матрицы представляет состояние автомата, а каждый столбец соответствует возможному входному элементу.
- ❑ В ячейках матрицы указывается структура из трех полей:
  - состояние
  - функция которая выполняется при переходе из одного состояния в другое
  - сообщение об ошибке

# 3.5 Матрица переходов КА

<десятичное целое с фиксированной точкой> ::=

<число с точкой> <цифра> |

< десятичное целое с  
фиксированной точкой > <цифра>

<число с точкой> ::= <целое>

<десятичная точка>

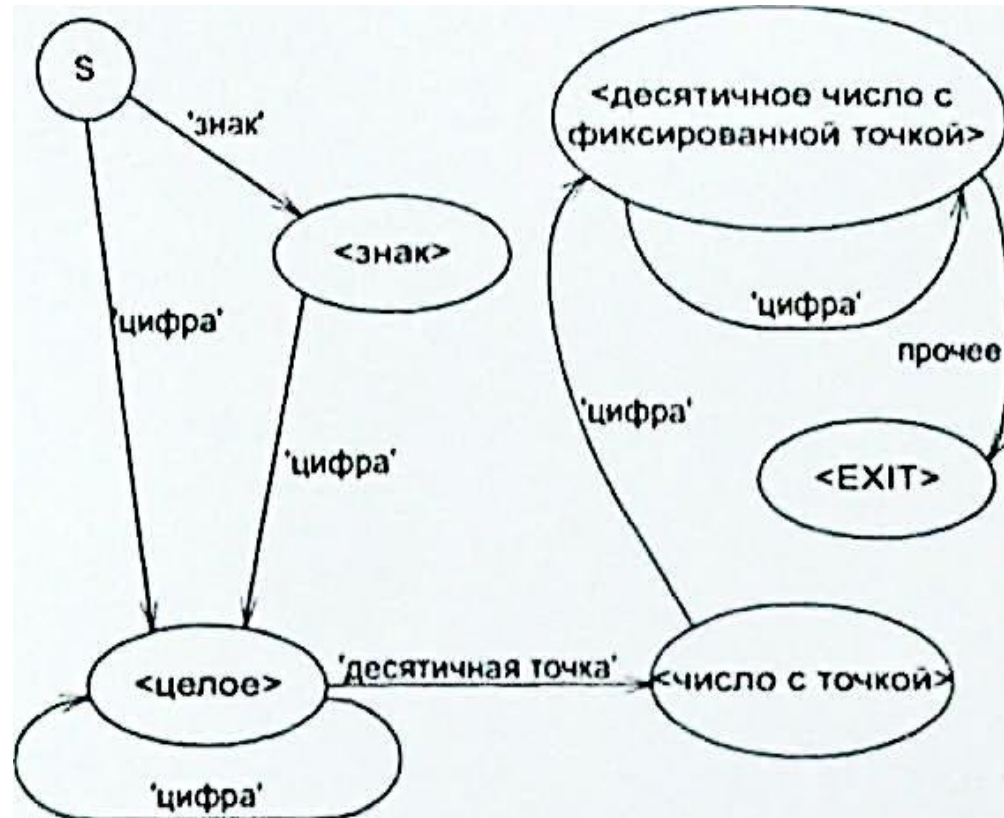
<целое> ::= <знак> <целое> |

<цифра> | <целое> <цифра>

<десятичная точка> ::= .

<знак> ::= + | -

<цифра> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9



## 3.5 Матрица переходов КА

Состояние\Следующий элемент	<знак>	<десятичная точка>	<цифра>	<прочие>
1 Start	2, F1, M0	6, F0, M1	3, F2, M0	6, F0, M1
2 <знак>	6, F0, M2	6, F0, M3	3, F2, M0	6, F0, M4
3 <целое>	6, F0, M5	4, F4, M0	3, F2, M0	6, F0, M4
4 <число с точкой>	6, F0, M6	6, F0, M1	5, F3, M0	6, F0, M6
5 <дес. число с фикс. т.>	6, F0, M8	6, F0, M1	5, F3, M0	Exit
6 Error	Exit	Exit	Exit	Exit

*Матрица переходов состояний для распознавания десятичных чисел с фиксированной точкой*

## 3.5 Матрица переходов КА

Подпрограмма	Действие
Инициализация	<code>num=0; sign = "+"; count = 1.0</code>
F0	Exit
F1	<code>sign = insign</code>
F2	<code>num = 10 * num + innum</code>
F3	<code>count = count / 10; num = num + innum * count</code>
F4	

`innum` – значение следующей введенной цифры;

`insign` – значение знака во входной строке;

`sign` – значение знака распознаного числа;

`num` – значение распознаного числа;

`count` – количество цифр после точки.

## **3.5 Матрица переходов КА**

<b>Номер сообщения</b>	<b>Сообщение</b>
М0	
М1	Ошибка. Строка не число: нет ни цифры ни знака
М2	Ошибка. Два последовательных знака
М3	Ошибка. Перед точкой нет ни одной цифры
М4	Ошибка. За знаком не следует цифра
М5	Ошибка. Отсутствует точка
М6	Ошибка. За точкой нет ни одной цифры
М7	Ошибка. Две десятичные точки в числе
М8	Ошибка. Знак в середине числа
Exit	All OK



## 3.6 Детерминированный и недетерминированный автомат

- Конечный автомат  $(Q, V, \delta, q_0, F)$  называется **детерминированным конечным автоматом (ДКА)**, если в каждом из его состояний для любого входного символа функция переходов содержит не более одного состояния.  
Для  $a \in V, q \in Q, r \in Q$ , определена  $\delta(a, q) = \{r\}$
- В недетерминированном КА  $\delta(a, q) = \{r_1, r_2, \dots, r_n\}$  – означает, что из состояния  $q$  по входному символу  $a$  можно осуществить переход в любое из состояний  $r_i, i = 1, 2, \dots, n$ .
- Доказано, что для любого КА можно построить ДКА.

**Домашнее задание: Алгоритм преобразования произвольного КА к детерминированному виду.**

## **3.7 Лексический анализатор**

- Лексический анализатор – программа, которая используя набор сканеров, преобразует исходный текст программы во внутреннее представление и помещает определённую информацию в специальные таблицы.
- Типы лексем (сканеров): идентификаторы, служебные слова, литералы (или константы) и разделители.
- Каждой лексеме сопоставляется пара:  
(тип лексемы, указатель на инф. о ней)

# 3.7 Лексический анализатор

- $\langle \text{prog} \rangle ::= \text{PROGRAM } \langle \text{prog\_name} \rangle \text{ VAR } \langle \text{dec\_list} \rangle \text{ BEGIN } \langle \text{stmt\_list} \rangle \text{ END.}$
- $\langle \text{prog\_name} \rangle ::= \text{id}$
- $\langle \text{dec\_list} \rangle ::= \langle \text{dec} \rangle ; | \langle \text{dec\_list} \rangle ; \langle \text{dec} \rangle$
- $\langle \text{dec} \rangle ::= \langle \text{id\_list} \rangle : \langle \text{type} \rangle$
- $\langle \text{type} \rangle ::= \text{INTEGER}$
- $\langle \text{id\_list} \rangle ::= \langle \text{id} \rangle | \langle \text{id\_list} \rangle , \langle \text{id} \rangle$
- $\langle \text{stmt\_list} \rangle ::= \langle \text{stmt} \rangle | \langle \text{stmt\_list} \rangle ; \langle \text{stmt} \rangle$
- $\langle \text{stmt} \rangle ::= \langle \text{assign} \rangle | \langle \text{read} \rangle | \langle \text{write} \rangle | \langle \text{for} \rangle$
- $\langle \text{assign} \rangle ::= \langle \text{id} \rangle := \langle \text{exp} \rangle$
- $\langle \text{exp} \rangle \langle \text{term} \rangle | \langle \text{exp} \rangle + \langle \text{term} \rangle | \langle \text{exp} \rangle - \langle \text{term} \rangle$
- $\langle \text{term} \rangle ::= \langle \text{factor} \rangle | \langle \text{term} \rangle * \langle \text{factor} \rangle | \langle \text{term} \rangle \text{ DIV } \langle \text{factor} \rangle$
- $\langle \text{factor} \rangle ::= \langle \text{id} \rangle | \langle \text{int} \rangle | ( \langle \text{exp} \rangle )$
- $\langle \text{read} \rangle ::= \text{READ}(\langle \text{id\_list} \rangle)$
- $\langle \text{write} \rangle ::= \text{WRITE}(\langle \text{id\_list} \rangle)$
- $\langle \text{for} \rangle ::= \text{FOR } \langle \text{index\_exp} \rangle \text{ DO } \langle \text{body} \rangle$
- $\langle \text{index\_exp} \rangle ::= \langle \text{id} \rangle := \langle \text{exp} \rangle \text{ TO } \langle \text{exp} \rangle$
- $\langle \text{body} \rangle ::= \langle \text{stmt} \rangle | \text{BEGIN } \langle \text{stmt\_list} \rangle \text{ END}$
- $\langle \text{id} \rangle ::= \langle \text{id\_name} \rangle$
- $\langle \text{id\_name} \rangle ::= \langle \text{letter} \rangle | \langle \text{id\_name} \rangle \langle \text{letter} \rangle | \langle \text{id\_name} \rangle \langle \text{digit} \rangle$
- $\langle \text{int} \rangle ::= \langle \text{digit} \rangle | \langle \text{int} \rangle \langle \text{digit} \rangle$
- $\langle \text{letter} \rangle ::= \text{A} | \text{B} | \text{C} | \dots | \text{Z}$
- $\langle \text{digit} \rangle ::= 0 | 1 | 2 | \dots | 9$

## 3.7 Лексический анализатор

```
PROGRAM STATS
VAR
    SUM, SUMSQ, I, VALUE, MEAN, VARIANCE: INTEGER;
BEGIN
    SUM := 0;
    SUMSQ := 0;
    FOR I 1 TO 100 DO
    BEGIN
        READ( VALUE );
        SUM := SUM + VALUE;
        SUMSQ := SUMSQ + VALUE * VALUE;
    END;
    MEAN := SUM DIV 100;
    VARIANCE := SUMSQ DIV 100 - MEAN * MEAN;
    WRITE( MEAN, VARIANCE);
END.
```

## 3.7 Лексический анализатор

Таблица служебных  
слов TW1

N	Терминальный символ	Признак разделителя
1	PROGRAM	-
2	BEGIN	-
3	VAR	-
4	INTEGER	-
5	END	-
6	FOR	-
7	TO	-
8	DO	-

Таблица разделителей  
TD 2

N	Терминальный символ	Признак разделителя
1	;	+
2	:	+
3	:=	

## 3.7 Лексический анализатор

Таблица литералов  
(констант) TNUM 3

N	Лите- рал	Ад- рес	Тип	Другая информа- ция
1	0			Значение
2	0			во
3	1			внутреннем
4	100			представлении

Таблица  
идентификаторов TID 4

N	Символи- ческое имя	Ад- рес	Тип	Другая информа- ция
1	STATS			
2	SUM			
3	SUMSQ			
4				

## **3.7 Лексический анализатор**

Переменные:

- `buf` – буфер для накопления символов лексем;
- `c` – очередной входной символ;
- `d` – переменная для формирования числового значения константы;
- `TW` – таблица служебных слов ;
- `TD` – таблица ограничителей;
- `TID` – таблица идентификаторов;
- `TNUM` – таблица констант;
- `tbl` – имя типа таблиц;
- `ptable` – указатель на `tbl`.

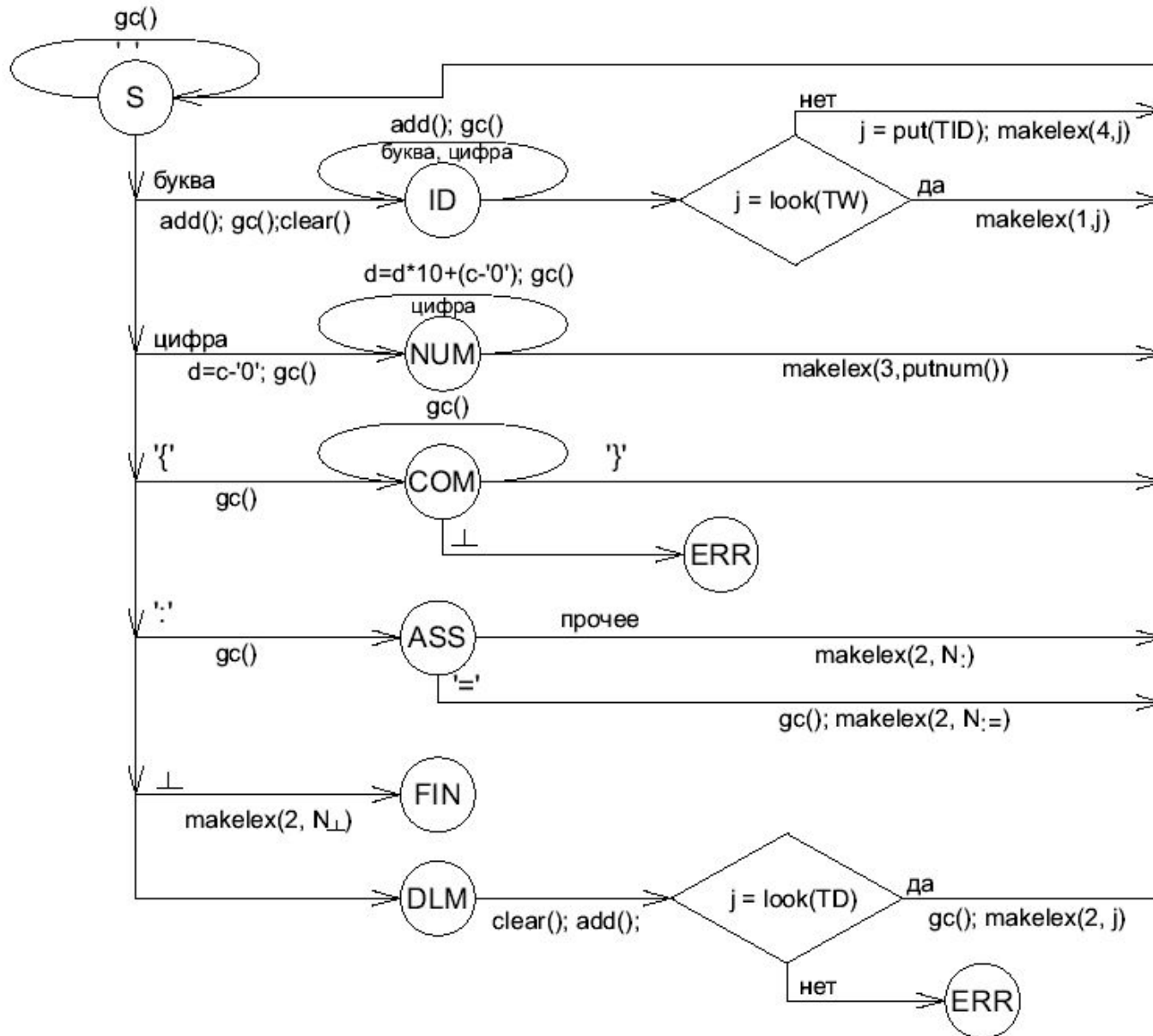
## 3.7 Лексический анализатор

Функции:

- `void clear (void)` – очистить буффер *buf*,
- `void add (void)` – добавление символа *c* в конец буфера *buf*;
- `int look (ptabl T)` – поиск в таблице *T* лексемы из буфера *buf*; результат – номер строки таблицы с информацией о лексеме либо 0, если лексемы нет в таблице;
- `int putl (ptabl T)` – запись в таблицу *T* лексемы из *buf*; результат – номер строки с информацией о лексеме;
- `int puntnum(void)` – запись в *TNUM* константы из *d* если ее там не было; результат – номер строки в таблице *TNUM* с информацией о константе-лексеме;
- `void makelex (int k, int i)` – формирование и вывод внутреннего представления лексемы; *k* – номер класса, *i* – номер в классе;
- `void gc (void)` – чтение из входного потока очередного символа и занесение его в *c*



# 3.7 Лексический анализатор



## 3.7 Лексический анализатор

- КА пользуясь набором сканера формирует набор лексем.

<1,1>      <4,1>      <1,3>      <4,2>

PROGRAM    STATS      VAR          SUM

<2,2>      <1,4>      <2,1>      <1,2>    ...

:          INTEGER    ;    BEGIN ...

## 3.7 Лексический анализатор

F0={ gc() }

F1={ clear(); add(); gc() }

F3={ add(); gc() }

F4={ j==look(TW) – ?

Да: makelex(1,j);

Нет: j=put(TID);

makelex(4,j); }

F5={ makelex(3, putnum()) }

F7={ makelex(2, look(TD)) }

F8={ add(),gc(), makelex(2,look(TD)) }

F9={ clear(), add()

j==look(TD) – ?

Да: gc(); makelex(2,j);

Нет: c=='\n' -?

Да:gc();

Нет:State=7; }

## 3.7 Лексический анализатор

		букв а	цифр а	' ,\n'	{	}	?	=	⊥	Проче е
0	S	1, F0	2, F1	3, F1	3, F0	7	4, F1	6	5	6
1	ID	1, F4	7	1, F3	0, F4	0, F4	0, F4	0, F4	0, F4	0, F4
2	NUM	0, F5	2, F3	0, F5	0, F5	0, F5	0, F5	0, F5	0, F5	0, F5
3	COM	3, F0	3, F0	3, F0	3, F0	0, F0	3, F0	3, F0	7	3, F0
4	ASS	0, F7	0, F7	0, F7	0, F7	0, F7	0, F7	0, F8	0, F7	0, F7
5	FIN	—	—	—	—	—	—	—	—	—
6	DLM	0, F9	0, F9	0, F9	0, F9	0, F9	0, F9	0, F9	0, F9	0, F9
7	ERR	—	—	—	—	—	—	—	—	—

## 3.8 Связь регулярных грамматик КА

- На основании имеющейся регулярной грамматики можно построить эквивалентный ей КА, и наоборот на основе КА можно построить эквивалентную ему грамматику.
- Поскольку КА являются распознавателями регулярных языков, таким образом, построив по регулярной грамматике КА, решаем задачу разбора.

### 3.8.1 Построение КА на основе леволинейной грамматики

Необходимо построить КА  $M(Q, V, \delta, q_0, F)$ , по леволинейной грамматике  $G(VT, VN, P, S)$ .

- 1) Множество входных символов строится на основании терминальных символов  $V=VT$
- 2) Множество состояний автомата строится на основании множества нетерминальных символов. Каждому нетерминалу ставится в соответствие состояние автомата и добавляется еще одно начальное состояние.

$$Q = VN \cup \{H\}$$

## 3.8.1 Построение КА на основе леволинейной грамматики

3) Просматриваем все множество правил грамматики  $G$ .

- Если встречается правило вида  $A \rightarrow t$ ,  $A \in VN$ ,  $t \in VT$ , то функцию перехода добавляем состояние  $A$ .  
 $A \in \delta(H, t)$ .
- Если  $A \rightarrow Bt$ ,  $A, B \in VN$ ,  $t \in VT$ , то  $A \in \delta(B, t)$ .

4) Начальное состояние автомата  $q_0 = H$ .

5) Множество конечных состояний включает одно состояние, соответствующее начальному символу грамматики  $S$ .

$$F = \{S\}$$

## 3.8.1 Построение КА на основе леволинейной грамматики

Построить автомат  $M(Q, V, \delta, q_0, F)$  на основе имеющейся автоматной грамматики:

$$G = ( \{0,1\}, \{S,A,B\}, P, S)$$

$$P: S \rightarrow A0 \mid B1$$

$$A \rightarrow S1 \mid 1$$

$$B \rightarrow S0 \mid 0$$

$$L(G) = \{ (10 \mid 01)^n, n > 0 \}$$

- $S \Rightarrow A0 \Rightarrow 10$
- $S \Rightarrow A0 \Rightarrow S10 \Rightarrow B110 \Rightarrow 0110$
- $S \Rightarrow B1 \Rightarrow 01$
- $S \Rightarrow B1 \Rightarrow S01 \Rightarrow A001 \Rightarrow 1001$



# 3.8.1 Построение КА на основе леволинейной грамматики

Шаг 1.  $V = \{0, 1\}$

Шаг 2.  $Q = \{S, A, B, H\}$

Шаг 3.  $\delta(A,0) = \{S\}$

$\delta(B,1) = \{S\}$

$\delta(S,1) = \{A\}$

$\delta(H,1) = \{A\}$

$\delta(S,0) = \{B\}$

$\delta(H,0) = \{B\}$

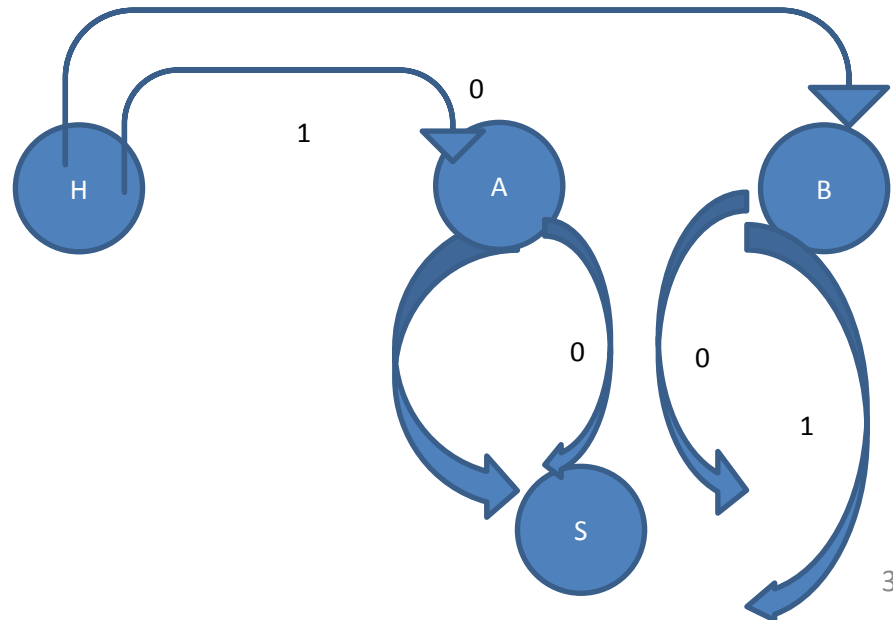
Шаг 4.  $q_0 = H$

Шаг 5.  $F = \{S\}$

$S \rightarrow A0 \mid B1$

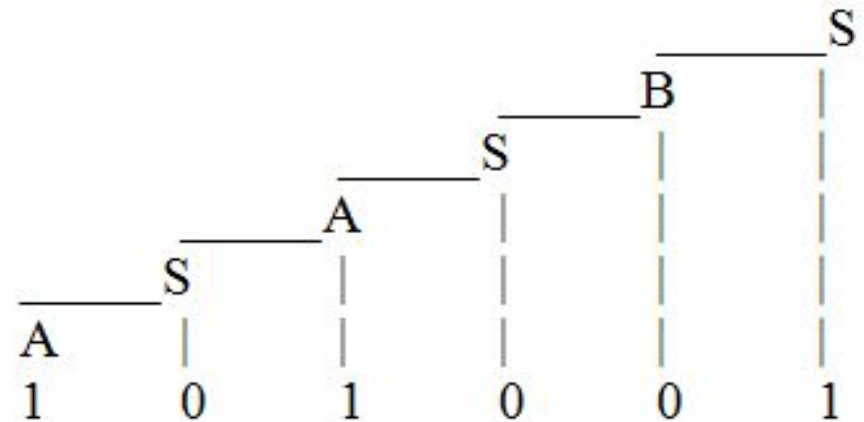
$A \rightarrow S1 \mid 1$

$B \rightarrow S0 \mid 0$



# 3.8.1 Построение КА на основе леволинейной грамматики

Шаг	Состояние	Остаток входной цепочки
1	H	101001
2	A	01001
3	S	1001
4	A	001
5	S	01
6	B	1
7	S	-



## **3.8.2 Построение леволинейной грамматики на основе КА**

1) Множество терминальных символов строится на основании входных символов

$$V_T = V$$

2) Множество нетерминалов грамматики  $G$  строится на основании множества состояний  $Q$  автомата за исключением начального состояния автомата

$$V_N = Q / \{q_0\}$$

## 3.8.2 Построение леволинейной грамматики на основе КА

3) Просматриваем функции переходов автомата  $M$  для всех возможных состояний из  $Q$  для всех входных символов из  $V$ .

Если  $\delta(A,t)=\{V_1, V_2, \dots, V_n\}$ ,  $A, V \in Q$ ,  $t \in V$ ,  $1 < i \leq n$ ,

то для всех состояний  $V_i$  выполняем следующее:

- Если  $A = q_0$ , то множество  $V_i \rightarrow t$
- Если  $A \neq q_0$ , то множество  $V_i \rightarrow A t$

4)

– Если множество конечных состояний  $F$  автомата  $M$  содержит только одно состояние  $F = \{f_0\}$ , то начальным символом грамматики принимается нетерминал, соответствующий этому состоянию.

– Если множество заключительных состояний

$$F = \{f_1, f_2, \dots, f_n\},$$

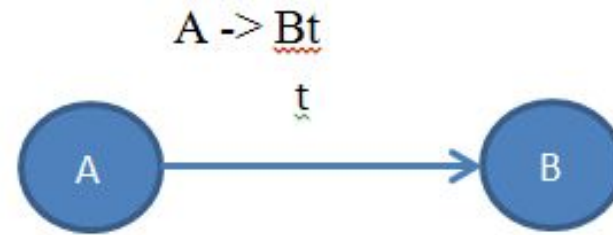
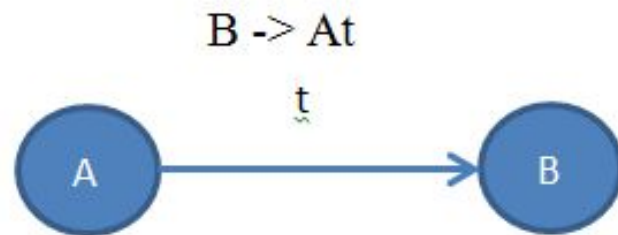
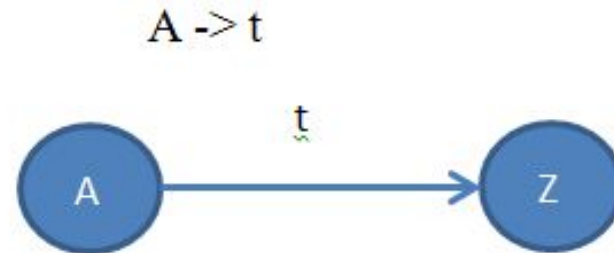
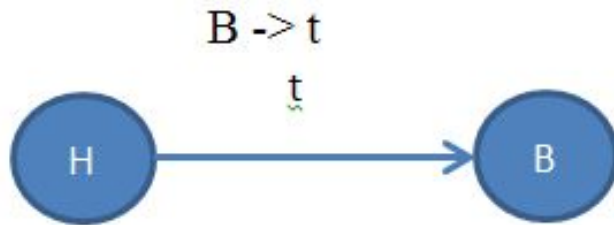
$$\text{то } S \rightarrow F_1 \mid F_2 \mid \dots \mid F_n$$

$$VN = VN \cup \{S\}$$

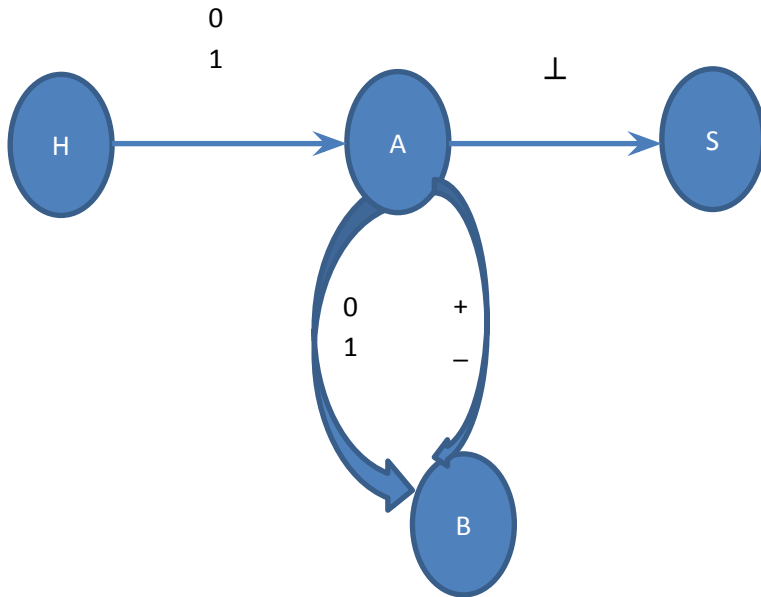
## 3.8.2 Построение леволинейной грамматики на основе КА

Для леволинейной грамматики

Для праволинейной грамматики



## 3.8.2 Построение леволинейной грамматики на основе КА



Шаг 1.  $V_T = V = \{0, 1, +, -, \perp\}$

Шаг 2.  $V_N = \{S, A, B\}$

Шаг 3.

$P: S \rightarrow A \perp$

$A \rightarrow A0 \mid A1 \mid B0 \mid B1 \mid 0 \mid 1$

$B \rightarrow A+ \mid A-$

Шаг 4.  $S = S$

$V = \{0, 1, +, -, \perp\}$

$Q = \{S, A, B, H\}$

$q_0 = H$

$F = \{S\}$

$\delta(B, 1) = \{A\}$

$\delta(A, 0) = \{A\}$

$\delta(B, 0) = \{A\}$

$\delta(A, 1) = \{A\}$

$\delta(A, +) = \{B\}$

$\delta(H, 0) = \{A\}$

$\delta(A, \perp) = \{S\}$

$\delta(A, -) = \{B\}$

$\delta(H, 1) = \{A\}$