

ЛЕКЦИЯ 5

“ПОСТРОЕНИЕ ЛЕКСИЧЕСКОГО АНАЛИЗАТОРА (СКАНЕРА) НА БАЗЕ КОНЕЧНОГО АВТОМАТА”

ПЛАН

1. Сканер как неотъемлемая часть компилятора
2. Процедура построения лексического анализатора (сканера) на базе конечного автомата.
3. Демонстрационный пример.
4. Синтаксические машины Глени.

Лексический анализатор (сканера)

Сканер – составная часть транслятора, производящая лексический анализ входной (исходной программы).

Лексический анализ программы – анализ синтаксиса на уровне лексических единиц или синтаксический анализ на уровне лексических единиц.

От английского scan [skæn] – глагол, обозначающий внимательно рассматривать, изучать.

В задачи сканера входят.

1. **Из цепочек литер** исходной программы **выделить** (распознать, составить) символы языка программирования, которые называются **атомами** или **лексемами** языка.

2. **Заменить последовательность литер** переменной длины **кодами фиксированной разрядности**, соответствующим опознанным лексемам.

3. **Готовить и сопровождать таблицы соответствия** объектов исходной программы: идентификаторов и констант системе кодирования лексем, номеров строк программы для облегчения диагностики.

4. **Представить** исходную **программу** в виде последовательности **кодов лексем (дескрипторов)**.

Считается, что лексема (или дескриптор лексем) состоит из двух частей: **класса** и **значения**.

Код класса
Код значения или ссылка

Процедура построение сканера

Сканер строится на основе теории интерпретирующих автоматов. Построение осуществляется в следующей традиционной последовательности.

1. **Определяется список лексем.**
2. **Определяются внутренние коды** лексем. Это могут быть как коды, используемые процессором, для которого создаётся транслятор, так и коды, произвольно назначенные программистом.
3. **Для лексем** анализируемого языка необходимо **построить описание** пригодное для построения конечного автомата.
4. **Строится обобщённая диаграмма** состояний сканера.
5. **Разрабатываются алгоритмы обработки** и процедуры их реализации, которые будут вызываться из различных состояний конечного автомата или при переходах между ними. Указанная информация наносится на обобщённую диаграмму.
6. **Разрабатываются тесты** для проверки корректности функционирования отдельных модулей и конечного автомата в целом. Определяются **диагностические сообщения** об ошибках.
7. **Составляются спецификации** программных модулей, входящих в сканер, и **осуществляется кодирование** алгоритмов на выбранном или заданном языке программирования.

Пример.

Предложения языка учебного микропроцессора описывается формальной грамматикой.

$G[Stmt]$:

$\langle Stmt \rangle$	$::=$	$\langle Lbl \rangle \langle Cop \rangle \langle Operand \rangle \langle Dir \rangle \langle Operand \rangle \langle Cop \rangle \langle Opernd \rangle $ $\langle Lbl \rangle \langle Cop \rangle \langle Cop \rangle \langle Lbl \rangle \langle Dir \rangle \langle Operand \rangle \langle Dir \rangle$
$\langle Cop \rangle$	$::=$	$add neg jmp shift$
$\langle Dir \rangle$	$::=$	$start org end$
$\langle Operand \rangle$	$::=$	$\# \langle data \rangle @ \langle iden \rangle @ \# \langle data \rangle$
$\langle iden \rangle$	$::=$	$\langle Bukva \rangle \langle Next \rangle \langle Bukva \rangle$
$\langle Bukva \rangle$	$::=$	$a b c \dots z$
$\langle Next \rangle$	$::=$	$\langle Bukva \rangle \langle Next \rangle \langle Cifra \rangle \langle Next \rangle \langle Bukva \rangle \langle Cifra \rangle$
$\langle Cifra \rangle$	$::=$	$0 1 2 \dots 9$
$\langle data \rangle$	$::=$	$0 \langle Bin \rangle b 0 \langle Hex \rangle h$
$\langle Bin \rangle$	$::=$	$\langle Bn \rangle \langle dig2 \rangle \langle dig2 \rangle$
$\langle dig2 \rangle$	$::=$	$0 1$
$\langle Hex \rangle$	$::=$	$\langle Hex \rangle \langle dig16 \rangle \langle dig16 \rangle$
$\langle dig16 \rangle$	$::=$	$0 1 2 \dots 9 a b c d e f$
$\langle Lbl \rangle$	$::=$	$\langle iden \rangle$

Список лексем и их коды:

Служебные слова			
Операции		Директивы	
<i>add</i>	101	<i>start</i>	201
<i>neg</i>	102	<i>org</i>	202
<i>jmp</i>	103	<i>end</i>	203
<i>shift</i>	104		
Разделители			
Однолитерные		Двулитерные	
#	301	@#	401
@	302		
Условно-терминальные символы			
Константа		Идентификатор	
<data>	501	<iden>	601

Общая структура регулярного выражения примерно такова:

$$Cc_1L_1 \cup Cc_2L_1 \cup \dots \cup Cc_iL_1 \cup \dots \cup Cc_rL_1 \cup \langle iden \rangle L_1 \cup \langle data \rangle L_2 \cup d_1 \cup d_2 \cup \dots \cup d_j \cup \dots \cup d_n \cup L_3.$$

В записи обозначены компоненты:

Cc_i – регулярное выражение для описания i -того служебного слова;

$\langle iden \rangle$ – регулярное выражение для описания классов переменных и меток;

$\langle data \rangle$ – регулярное выражение для описания класса констант;

L_1 – лексема, состоящая в том, что текущая литера не буква или не цифра;

L_2 – текущая литера не цифра;

d_j – разделители j -того типа;

L_3 – литера, не принадлежащая алфавиту конечного автомата.

Переменную $\langle iden \rangle$ в эскизном виде можно представить регулярным выражением

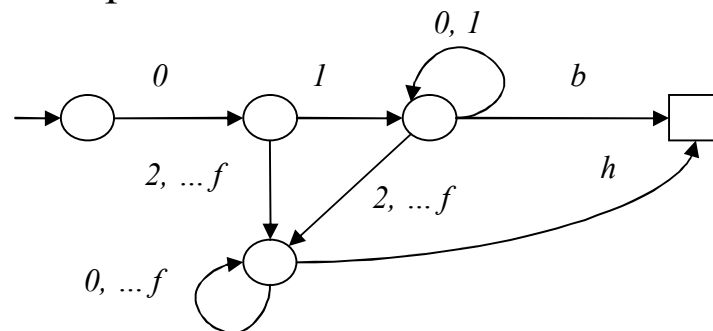
$$B \{B \vee C\},$$

где B и C обозначают букву латинского алфавита и арабскую цифру соответственно.

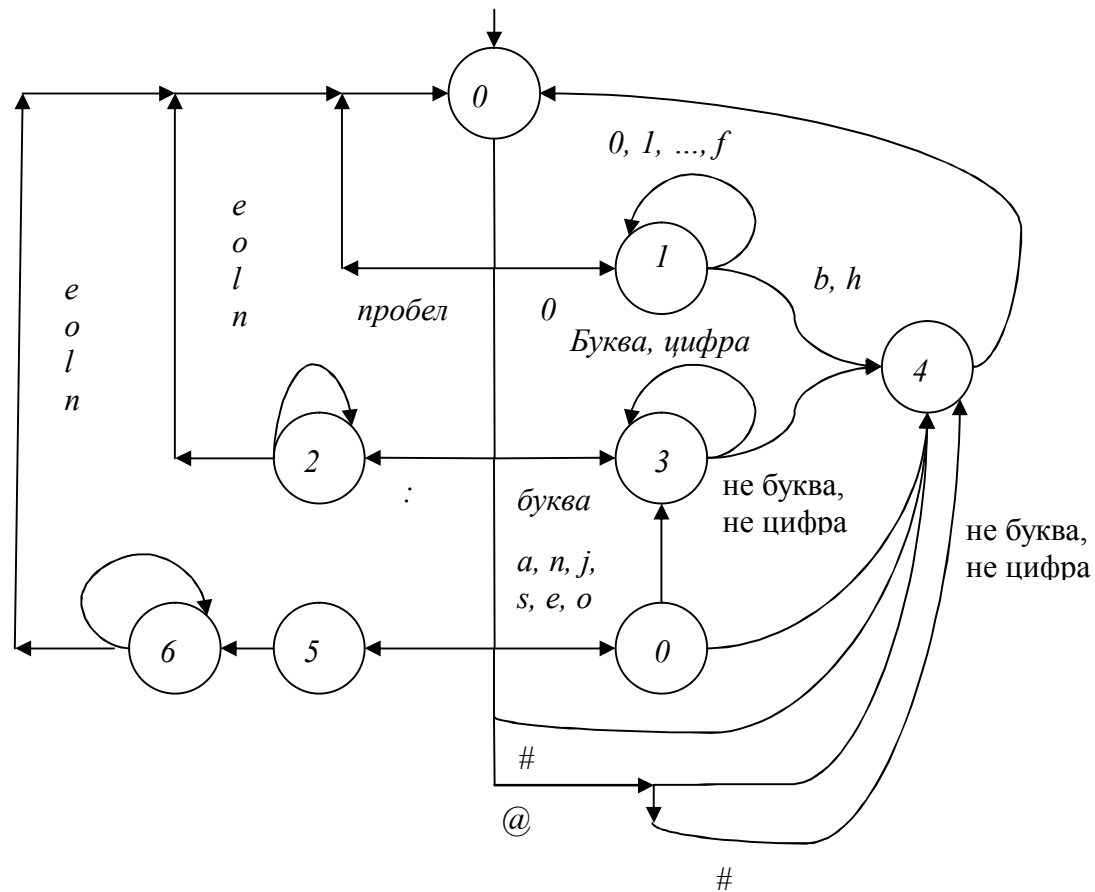
Арифметические константы $\langle data \rangle$ описываются так

$$0\{0 \vee 1\}b \vee 0\{0 \vee 1 \vee 2 \vee 3 \vee 4 \vee 5 \vee 6 \vee 7 \vee 8 \vee 9 \vee a \vee b \vee c \vee d \vee e \vee f\}h.$$

Фрагмент КА для анализа чисел

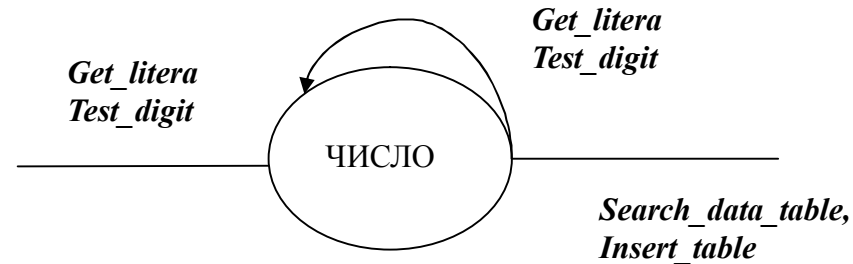


Общая схема сканера



Обозначено: 0 – начальное состояние; 1 – распознавание константы; 2 – игнорирование комментария; 3 – распознавание идентификатора; 4 – формирование лексем и ведение таблиц; 5 – ошибка, встречен символ, который не принадлежит алфавиту; 6 – игнорирование оператора, в котором была допущена ошибка; 7 – распознавание служебных слов.

Нанесение на диаграмму процедур поддержки и обработки приведёт примерно к следующей картине для каждого состояния конечного автомата



Предполагается разработка следующих процедур: ***Get_litera*** – чтение литеры из потока; ***Test_digit*** – проверка на принадлежность литеры к классу цифр; ***Search_data_table*** – определение, является ли сформированный символ новым или уже присутствует в таблице; ***Insert_Table*** – запись нового символа в таблицу.

РАЗРАБОТКА ТЕСТОВОГО ПРИМЕРА

На вход сканера поступает программа на учебном ассемблере.

```

        org      @#030eh ; стартовый адрес
start
...
m1      add      @alpha
...
end

```

В результате ожидается, что будут построены таблицы констант и идентификаторов

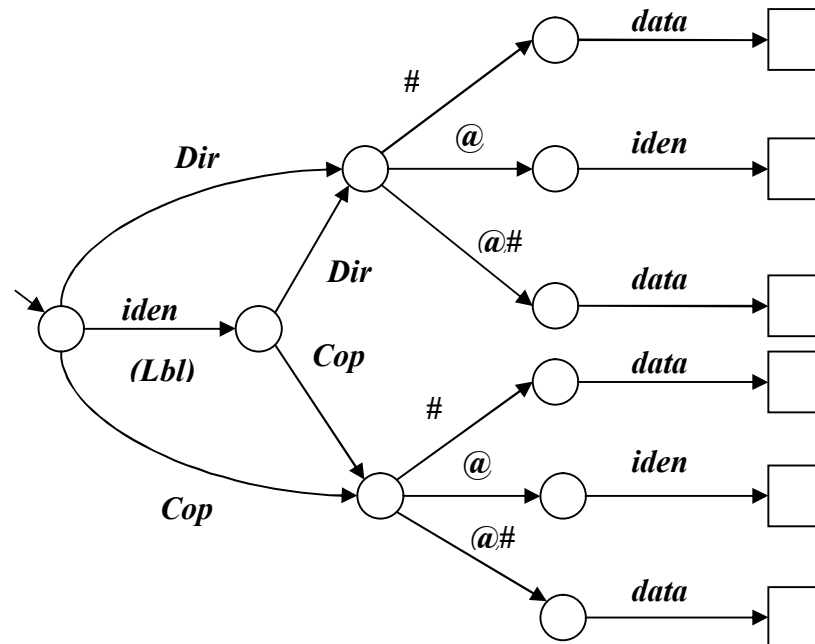
Таблица констант		
№	Константа	Значение
1	30e	0000 0011 0000 1110
...
42	157	0001 0101 0111

Таблица имён	
№	Имя переменной
1	x
...	...
27	alpha
...	...
48	m1

После перекодировки должна получиться следующая последовательность лексем

PROG.COD	202	401	501	...	601	101	302	601	...	203
PROG.NUM	000	000	001	...	048	000	000	049	...	000
ЛЕКCEMA	org	@#	030eh		m1	add	@	alpha		end

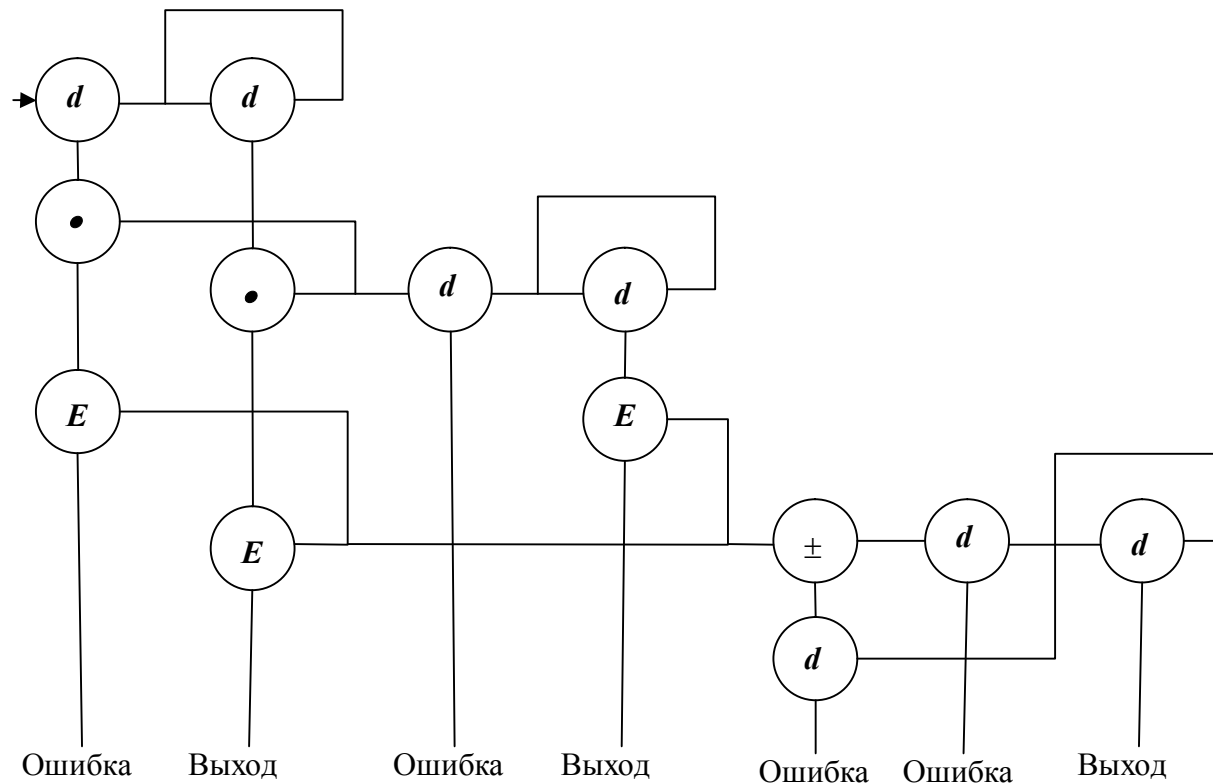
КОНЕЧНЫЙ АВТОМАТ ДЛЯ СИНТАКСИЧЕСКОГО АНАЛИЗА



СИНТАКСИЧЕСКИЕ МАШИНЫ ГЛЕННИ

Glennie A.E. в 1960 г. Кругом обозначено сравнение, выход направо – положительный результат сравнения, при этом входная строка продвигается на одну позицию (1 литеру), выход вниз – отрицательный, продвижение входной строки блокируется

Неминимальная машина Гленни



Машина Гленни с минимальным числом состояний

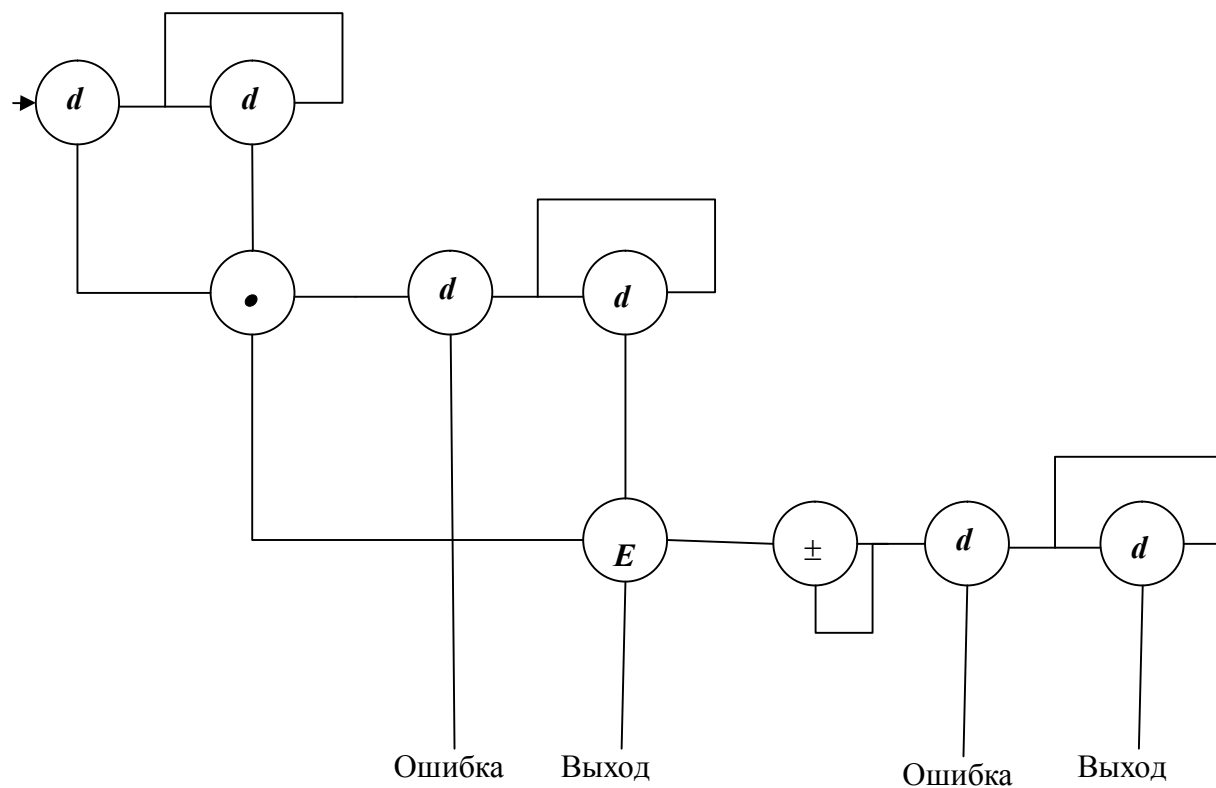


Таблица управления минимальной синтаксической машины

Номер состояния	Символ на входе	Следующее состояние		Выполняемая операция	
		True	False	True	False
1	<i>d</i>	2	3	V	—
2	<i>d</i>	2	3	V	—
3	•	4	6	—	—
4	<i>d</i>	5	Error	W	Error
5	<i>d</i>	5	6	W	—
6	<i>E</i>	7	Exit	—	Y
7	+ √ -	8	8	Z	—
8	<i>d</i>	9	Error	X	Error
9	<i>d</i>		Exit	X	Y

Фрагменты алгоритмов процедур обработки

```

Procedure V /* подсчёт мантиссы целой части числа при вводе
/* сдвиг организован умножением на 10 */
Mantissa := 10 * Mantissa + d
End V

```

```

Procedure W /* подсчёт мантиссы дробной части числа при вводе */
N := N + 1 /* подсчёт длины дробной части */

```

```

Mantissa := 10 * Mantissa + d
End W
  Procedure X
    /* подсчёт показателя степени в научной нотации числа */
    Expo := 10 * Expo + d
  End X
Procedure Y
/* окончательное получение эквивалента числа в формате с плавающей
точкой */
Result := Mantissa * 10 ^ (Znak * Expo - N)
End Y
  Procedure Z
    /* определение знака в показателе степени числа в научной
нотации */
    Znak := sign
  End Z
Procedure Error /* выдача диагностического сообщения */
output ("Ошибка в записи числа")
End Error

```

Перед началом работы переменные *Mantissa*, *N*, *Expo* положены равными нулю, а переменная *Znak* равна единице, что соответствует знаку плюсу в показателе степени.