

Лабораторная работа №3

Регулярные грамматики и конечные автоматы

Цель работы состоит в практическом построении лексического анализатора, распознающего цепочки символов генерируемых регулярными грамматиками.

Краткие теоретические сведения

Фазы компиляции включают лексический анализ, синтаксический анализ, семантический анализ, фазу генерации внутреннего представления программы, оптимизацию кода и генерацию объектного модуля.

В реальных компиляторах порядок фаз компиляции может быть несколько иным. Некоторые из них, например, лексический анализ и синтаксический анализ, целесообразно объединить в одну фазу. В интерпретаторах некоторые этапы компиляции могут отсутствовать.

В основе построения лексических анализаторов лежат регулярные грамматики. В настоящей лабораторной работе под регулярной грамматикой понимается *леволинейная* грамматика, не принадлежащая иерархии Хомского¹. Для грамматик этого типа существует алгоритм разбора, который отвечает на вопрос: принадлежит ли анализируемая цепочка языку, порождаемому этой грамматикой, или нет. Алгоритм разбора таков:

- ✓ первый символ исходной цепочки $w = a_1 a_2 \dots a_n$ заменяется нетерминалом A , для которого в грамматике есть правило вывода $A \rightarrow a_1$ (производится «свертка» терминала a_1 к нетерминалу A);
- ✓ до тех пор, пока не будет прочтена вся цепочка, выполняются следующие шаги:
- ✓ если в правилах грамматики имеется продукция $B \rightarrow A a_i$ ($i=2, 3, \dots, n$), то полученный на предыдущем шаге разбора нетерминал A , а также расположенный непосредственно справа от него очередной терминал a_i входной цепочки заменяется нетерминалом B .

Это эквивалентно построению дерева разбора восходящим методом: на каждом шаге алгоритма строится один из уровней дерева разбора от листьев к корню.

При работе этого алгоритма возможны следующие ситуации:

¹ Леволинейные грамматики не используются на практике.

- прочитана вся цепочка; на каждом шаге находилась единственная нужная «свертка»; на последнем шаге свертка привела к символу S . Это означает, что исходная цепочка $w = a_1 a_2 \dots a_n \perp$ принадлежит языку $L(\mathcal{G})$ генерируемому грамматикой \mathcal{G} .
- прочитана вся цепочка; на каждом шаге находилась единственная нужная «свертка»; на последнем шаге «свертка» не привела к символу S . Это означает, что исходная цепочка отвергается, она не принадлежит языку $L(\mathcal{G})$.
- на некотором шаге не нашлось нужной «свертки», т.е. для полученного на предыдущем шаге нетерминала A и расположенного непосредственно справа от него очередного терминала a_i исходной цепочки не нашлось нетерминала B , для которого в грамматике было бы правило вывода $B \rightarrow A a_i$. Это означает, что исходная цепочка отвергается, она не принадлежит языку $L(\mathcal{G})$.
- на некотором шаге работы алгоритма оказалось, что есть более одной подходящей «свертки», т.е. в грамматике разные нетерминалы имеют правила вывода с одинаковыми правыми частями, и поэтому непонятно, к какому из них производить «свертку». Это говорит о *недетерминированности разбора*. Анализ этой ситуации будет дан ниже.

Пусть разбор на каждом шаге детерминированный. Распознаватель полностью характеризуется таблицей, строки которой помечены нетерминальными символами грамматики, а столбцы – терминальными. Значение каждого элемента таблицы – это нетерминальный символ, к которому можно свернуть конкатенацию нетерминала и терминала, помечающие соответствующие строку и столбец.

Пример 1. Для грамматики $\mathcal{G} = (\{a, b, \perp\}, \{S, A, B, C\}, P, S)$ с правилами

$P = \{S \rightarrow C \perp, C \rightarrow Ab \mid Ba, B \rightarrow b \mid Cb\}$ таблица переходов такова

Табл. Таблица переходов

	a	b	\perp
C	A	B	S
A	-	C	-
B	C	-	-
S	-	-	-

Прочерк ставится в том случае, если «свертки» нет. Таблице переходов в соответствие ставится направленный помеченный графу, который показан на рис. 1.

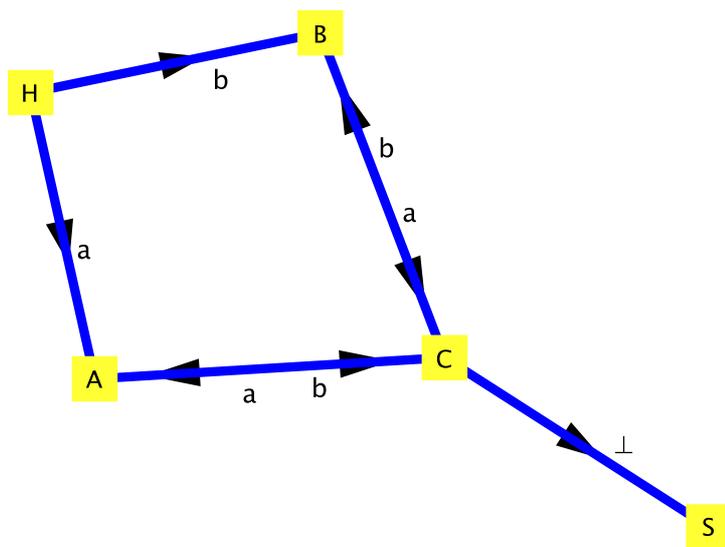


Рис. 1: Диаграмма переходов из **Примера 1**

Граф строится следующим образом:

✓ Строятся вершины графа, помеченные нетерминалами грамматики. Для каждого нетерминала грамматики граф содержит только одну вершину. Добавляется еще одна вершина, помеченная символом, отличным от символов нетерминалов, например, H . Эти вершины называют состояниями, H – начальное состояние.

✓ Вершины соединяются дугами по следующим правилам:

- для каждого правила грамматики вида $A \rightarrow a$ дугой соединяются состояния от H к A . Дуга помечается символом a .
- для каждого правила $A \rightarrow Ba$ дугой соединяются состояния от B к A , дуга помечается символом a .

Цепочки языка порождаются выводом в грамматике, например, для грамматики \mathcal{G} с правилами $\mathbf{P} = \{S \rightarrow C \perp, C \rightarrow Ab \mid Va, B \rightarrow b \mid Cb, A \rightarrow a \mid Ca\}$ из **Примера 1** вывод таков $S \Rightarrow C \perp \Rightarrow Ab \perp \Rightarrow Cab \perp \Rightarrow Vaab \perp \Rightarrow Cbaab \perp \Rightarrow Vabaab \perp \Rightarrow babaab \perp$.

Пусть имеются леворекурсивные правила $A \rightarrow A\gamma \mid \alpha$. Цепочки символов, порождаемые этими двумя правилами таковы $w = \alpha\gamma^* \in \mathbf{L}(\mathcal{G})$. Эти же цепочки можно порождаются другими праволинейными правилами $A \rightarrow \alpha A' \mid \gamma A' \mid \varepsilon$, где A' – дополнительный нетерминал.

Алгоритм разбора по диаграмме состояний:

1. Начальным текущим состоянием объявляется H .

2. Пока цепочка полностью не прочитана, выполняются следующие шаги. Считывается текущий символ входной цепочки и осуществляется переход из текущего состояния в другое состояние по дуге, помеченной этим символом цепочки. Новое состояние становится текущим.
3. При работе этого алгоритма возможны следующие исходы:
 - I. Прочитана вся цепочка. На каждом шаге находилась единственная дуга, помеченная очередным символом анализируемой цепочки. В результате последнего перехода распознаватель пришел в состояние S . Это означает, что исходная цепочка принадлежит языку $L(\mathcal{G})$.
 - II. Прочитана вся цепочка. Распознаватель пришел в состояние, отличное от состояния S . Это означает, что исходная цепочка не принадлежит языку $L(\mathcal{G})$.
 - III. На некотором шаге не нашлось дуги, выходящей из текущего состояния и помеченной очередным анализируемым символом. Это означает, что исходная цепочка не принадлежит языку $L(\mathcal{G})$.
 - IV. На некотором шаге работы алгоритма оказалось, что есть несколько дуг, выходящих из текущего состояния, помеченных очередным анализируемым символом, но ведущих в разные состояния. Это говорит о *недетерминированности разбора*. Анализ этой ситуации будет приведен ниже.

Диаграмма состояний определяет конечный автомат, построенный по регулярной грамматике, который допускает множество цепочек, генерируемых грамматикой. Состояния и дуги ДС – это графическое изображение функции переходов конечного автомата из состояния в состояние. Очередной анализируемый символ совпадает с символом-меткой дуги перехода.

Среди всех состояний выделяется начальное состояние. Вначале автомат находится в этом состоянии. Всегда имеется конечное состояние, в котором автомат завершает работу. Переходом в это состояние автомат допускает проанализированную им цепочку.

Определение. Конечный автомат (КА) – это пятерка $(Q, \Sigma, \delta, F, S)$, где Q – конечное множество состояний; Σ – конечное множество допустимых входных символов; δ – отображение множества $Q \times \Sigma$ во множество Q , то есть $\delta: Q \times \Sigma \rightarrow Q$; $S \in Q$ – начальное состояние КА; $F \subset Q$ – конечное множество заключительных состояний. Значение функции перехода $\delta(A, a) = B$ означает, что из состояния $A \in Q$ по входному символу $a \in \Sigma$ происходит переход в состояние B . Запись $\delta(A, a) = \emptyset$ говорит, что следующий шаг автомата невозможен.

Конечный автомат *допускает цепочку* $a_1 a_2 \dots a_n$, если $\delta(H, a_1) = A_1$, $\delta(A_1, a_2) = A_2$, ..., $\delta(A_{n-2}, a_{n-1}) = A_{n-1}$, $\delta(A_{n-1}, a_n) = S$, где $a_i \in \Sigma$, $A_j \in Q$ ($j=1, 2, \dots, n-1; i=1, 2, \dots, n$), H – начальное состояние, S – одно из заключительных состояний.

Для работы с диаграммами состояний вводятся соглашения:

- ✓ если из одного состояния в другое выходит несколько дуг, помеченных разными символами, то будем изображать одну дугу, помеченную всеми этими символами;
- ✓ непомеченная дуга будет соответствовать переходу при любом символе, кроме тех, которыми помечены другие дуги, выходящие из этого состояния.
- ✓ введем состояние ошибки (ER); переход в это состояние будет означать, что исходная цепочка языку не принадлежит.

По диаграмме состояний написан анализатор для регулярной грамматики. Для грамматики из **Примера 1** анализатор будет таким:

```
#include <stdio.h>
#include <iostream>
int main()
{
    std::cout << "Labo #3. Finite automata recognizer.\nGrammar with left-recursive
productions: \nS->Cc\nC->Ab|Ba\nA->Ba\nB->Ab" << std::endl;
    enum states {H, A, B, C, S, ER}; /* the set of states */
    enum states CS; /* CS is a current state */
    FILE *fp; /*file pointer */
    char c;
    int cnt = 0;
    std::string name;
    name = "data.txt";
    CS = H;
    std::cout << "Data file: " << name << "\n\n\n";
    fp = fopen("data.txt","r");
    c = fgetc(fp);
    cnt++;
    printf("char %c number %d\n", c,cnt);
    do
    {
        switch(CS)
        {
            case H:
                if(c == 'a')
                {
                    c = fgetc(fp);
                    cnt++;
                    printf("char %c number %d\n", c,cnt);
                    CS = A;
                }
                else if(c == 'b')
                {
                    c = fgetc(fp);
                    cnt++;
                    printf("char %c number %d\n", c,cnt);
                    CS = B;
                }
                else CS = ER;
                break;
            case A:
                if(c == 'b')
                {
                    c = fgetc(fp);
                    cnt++;
                    printf("char %c number %d\n", c,cnt);
                    CS = C;
                }
                else CS = ER;
                break;
```

```

case B:
    if(c == 'a')
        {
            c = fgetc(fp);
            cnt++;
            printf("char %c number %d\n", c,cnt);
            CS = C;
        }
    else CS = ER;
    break;
case C:
    if(c == 'a')
        {
            c = fgetc(fp);
            cnt++;
            printf("char %c number %d\n", c,cnt);
            CS = A;
        }
    else if(c == 'b')
        {
            c = fgetc(fp);
            cnt++;
            printf("char %c number %d\n", c,cnt);
            CS = B;
        }
    else if(c == 'c') CS = S;
    else CS = ER;
    break;
}
} while(CS != S && CS != ER);
if(CS == ER)
    {
        printf("Rejected chain in position %d\n", cnt);
        return -1;
    }
else
    {
        printf("Accepted chain of length %d\n", cnt);
        return 0;
    }
}

```

При анализе по регулярной грамматике может оказаться, что несколько нетерминалов имеют одинаковые правые части, и поэтому неясно, к какому из них делать «свертку» (см. описание алгоритма). В терминах диаграммы состояний это означает, что из одного состояния выходит несколько дуг, ведущих в разные состояния, но помеченных одним и тем же символом.

Пример 2. Для грамматики $\mathcal{G} = (\{a, b, \perp\}, \{S, A, B\}, \mathbf{P}, S)$ с множеством продукций $\mathbf{P} = \{S \rightarrow Aa, A \rightarrow a \mid Bb, B \rightarrow b \mid Bb\}$, распознаватель будет недетерминированным, поскольку нетерминалы A и B обладают одинаковыми правыми частями – Bb . Такой распознаватель называется недетерминированным конечным автоматом.

Определение. Недетерминированный конечный автомат (НКА) – это пятерка $(\mathbf{Q}, \Sigma, \delta, \mathbf{F}, S)$, где \mathbf{Q} – конечное множество состояний; Σ – конечное множество допустимых входных символов; δ – отображение множества $\mathbf{Q} \times \Sigma$ во множество подмножеств \mathbf{Q} , то есть $\delta: \mathbf{Q} \times \Sigma \rightarrow \mathbb{P}(\mathbf{Q})$; $S \in \mathbf{Q}$ – начальное состояние НКА; $\mathbf{F} \subset \mathbf{Q}$ –

конечное множество заключительных состояний. Значение функции перехода $\delta(A, a) = \{B_1, B_2, \dots, B_n\}$ означает, что из состояния $A \in Q$ по входному символу $a \in \Sigma$ можно осуществить переход в любое из состояний $B_i \in Q$ ($i = 1, 2, \dots, n$), запись $\delta(A, a) = \emptyset$ говорит, что следующий шаг автомата невозможен.

На практике НКА не используются, поскольку моделируются с использованием переборных алгоритмов. НКА играют ключевую роль в теории формальных языков. В теории формальных языков доказана теорема, что для любого НКА всегда можно построить детерминированный КА, распознающий тот же язык, распознает и НКА.

Индивидуальные задания

1	<p>Дана грамматика: $G = (\{Q, P, R, S\}, \{0, 1, *, \perp, /\}, V, Q)$, где V: $S \rightarrow 0R$ $R \rightarrow /Q \mid \perp$ $Q \rightarrow 1P$ $P \rightarrow *S \mid \perp$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
2	<p>Дана грамматика: $S \rightarrow C\perp$ $B \rightarrow B1 \mid 0 \mid D0$ $C \rightarrow B1 \mid C1 \mid D \rightarrow D0 \mid 0$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
3	<p>Дана грамматика: $S \rightarrow C\perp$ $C \rightarrow B1$ $B \rightarrow 0 \mid D0$ $D \rightarrow B1$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
4	<p>Дана грамматика: $S \rightarrow A0$ $A \rightarrow A0 \mid S1 \mid 0$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
5	<p>Дана грамматика: $S \rightarrow 0A \mid 1S$ $A \rightarrow 0A \mid 1B$ $B \rightarrow 0B \mid 1B \mid \perp$ Определить язык, который она порождает; построить ДС;</p>

	написать на Си анализатор.
6	<p>Дана грамматика:</p> $S \rightarrow B \perp$ $A \rightarrow B1 0$ $B \rightarrow A1 C1 B0 1$ $C \rightarrow A0 B1$ <p>Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
7	<p>Дана грамматика:</p> $G = (\{S, C, D\}, \{0, 1\}, P, S), \text{ где } P:$ $S \rightarrow 1C 0D$ $C \rightarrow 0D 0S 1$ $D \rightarrow 1C 1S 0$ <p>Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
8	<p>Дана грамматика:</p> $S \rightarrow B0 0$ $B \rightarrow B0 C1 0 1$ $C \rightarrow B0$ <p>Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
9	<p>Дана грамматика:</p> $S \rightarrow Sb Aa a b$ $A \rightarrow Aa Sb a$ <p>Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
10	<p>Дана грамматика:</p> $S \rightarrow A0 A1 B1 0 1$ $A \rightarrow A1 B1 1$ $B \rightarrow A0$ <p>Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
11	<p>Дана грамматика:</p> $S \rightarrow S0 A1 0 1$ $A \rightarrow A1 B0 0 1$ $B \rightarrow A0$ <p>Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
12	<p>Дана грамматика:</p> $S \rightarrow A \perp$ $A \rightarrow A0 A1 B1$ $B \rightarrow B0 C0 0$ $C \rightarrow C1 1$ <p>Определить язык, который она порождает; построить ДС; написать на Си анализатор.</p>
13	<p>Дана грамматика:</p> $G = (\{I, J, K, M, N\}, \{0, 1, \sim, !\}, P, I), \text{ где } P:$ $I \rightarrow 0J 1K 0M$

	$J \rightarrow \sim K \mid 0M$ $K \rightarrow \sim M \mid 0J \mid 0N$ $M \rightarrow 1K \mid !$ $N \rightarrow 0! \mid 1! \mid !$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.
14	Дана грамматика: $G = (\{S, A, B, C\}, \{a, b, c\}, P, S)$, где P: $S \rightarrow aA \mid bB \mid aC$ $A \rightarrow bA \mid bB \mid c$ $B \rightarrow aA \mid cC \mid b$ $C \rightarrow bB \mid bC \mid a$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.
15	Дана грамматика: $G = (\{K, L, M, N\}, \{a, b, +, -, \perp\}, P, K)$, где P: $K \rightarrow aL \mid bM$ $L \rightarrow -N \mid -M$ $M \rightarrow +N$ $N \rightarrow aL \mid bM \mid \perp$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.
16	Дана грамматика: $G = (\{X, Y, Z, V, W\}, \{0, 1, x, y, z\}, P, X)$, где P: $X \rightarrow yY \mid zZ$ $Y \rightarrow 1V$ $Z \rightarrow 0W \mid 0Y$ $V \rightarrow xZ \mid xW \mid 1$ $W \rightarrow 1Y \mid 0$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.
17	Дана грамматика: $G = (\{S, A, B, C, D\}, \{a, b, c, d, \perp\}, P, S)$, где P: $S \rightarrow aA \mid bB$ $A \rightarrow cC \mid \perp$ $C \rightarrow cC \mid cA$ $B \rightarrow dD \mid \perp$ $D \rightarrow dD \mid dB$ Определить язык, который она порождает; построить ДС; написать на Си анализатор.
18	Дана грамматика: $G = (\{S, A, B, C, D\}, \{a, b, c, d, \#, \perp\}, P, S)$, где P: $S \rightarrow aA \mid bB \mid cC$ $A \rightarrow dD$ $B \rightarrow \#D$ $D \rightarrow dD \mid dB \mid \perp$ $C \rightarrow cD$ Определить язык, который она порождает; построить ДС;

	написать на Си анализатор.
19	<p>Даны две грамматики G1 и G2, порождающие языки L1 и L2. Построить регулярную грамматику для $L1 \cup L2$. Построить ДС и написать на Си анализатор для полученной грамматики.</p> <p>G1: $S \rightarrow S1 \mid A0$ $A \rightarrow A1 \mid 0$</p> <p>G2: $S \rightarrow A1 \mid B0 \mid E1$ $A \rightarrow S1$ $B \rightarrow C1 \mid D1$ $C \rightarrow 0$ $D \rightarrow B1$ $E \rightarrow E0 \mid 1$</p>
20	<p>Даны две грамматики G1 и G2, порождающие языки L1 и L2. Построить регулярную грамматику для пересечения L1 и L2. Для полученной грамматики построить ДС и написать на Си анализатор.</p> <p>G1: $S \rightarrow S1 \mid A0$ $A \rightarrow A1 \mid 0$</p> <p>G2: $S \rightarrow A1 \mid B0 \mid E1$ $A \rightarrow S1$ $B \rightarrow C1 \mid D1$ $C \rightarrow 0$ $D \rightarrow B1$ $E \rightarrow E0 \mid 1$</p>
21	<p>Для данной грамматики</p> <ul style="list-style-type: none"> • определить ее тип; • определить, какой язык она порождает; • написать регулярную грамматику, эквивалентную данной; • построить ДС и анализатор на Си. <p>$S \rightarrow 0S \mid S0 \mid D$</p> <p>$D \rightarrow DD \mid 1A \mid \varepsilon$</p> <p>$A \rightarrow 0B \mid \varepsilon$</p> <p>$B \rightarrow 0A \mid 0$</p>

Литература

- [1] Молчанов А.Ю. Системное программное обеспечение: Учебник для вузов. 3-е изд. – СПб.: Питер, 2010. – 400 с.
- [2] Ахо А.В., Лам М. С., Сети Р., Ульман Дж. Д. Компиляторы: принципы, технологии и инструментарий, 2-е изд.: Пер. с англ. – М.: Вильямс, 2008. – 1184 с. (разделы 3.4,3.6, 3.7)
- [3] Волкова И.А., Руденко Т.В. Формальные грамматики и языки. Элементы теории трансляции: Учебное пособие. – Издательский отдел факультета вычислительной математики и кибернетики МГУ им. М.В. Ломоносова, 1999. – 62 с.
- [4] Свердлов С.З. Языки программирования и методы трансляции: Учебное пособие. — СПб.: Питер, 2007. — 638 с. (стр. 221–230)