

А.С. СЕМЕНОВ

**ПРАКТИЧЕКИЙ КУРС ПО ТЕОРИИ АВТОМАТОВ И
ФОРМАЛЬНЫХ ЯЗЫКОВ**

С ПРИМЕРАМИ НА ЯЗЫКЕ ПРОГРАММИРОВАНИЯ C#

Москва 2019

СОДЕРЖАНИЕ

Введение	4
Глава 1. Классификация грамматик и автоматов	5
1.1. Способы описания синтаксиса языков	6
1.2. Классификация Хомского.	9
Глава 2. Автоматные грамматики и конечные автоматы	16
2.1. Определение автоматных грамматик и конечных автоматов ...	30
2.2. Свойства регулярных языков: лемма о накачке	30
2.3. Способы задания конечных автоматов	33
2.4. Пример проектирования грамматики и распознавателя	39
2.5. Способы реализации: по заданному регулярному выражению составные автоматы	52
2.6. Способы программной реализации	53
Глава 3. Контекстно-свободные грамматики и МП-автоматы	55
3.1. Определение КС-грамматик и МП-автоматов	60
3.2. Преобразование КС-грамматик	65
3.3. Пример преобразования КС-грамматик к приведенной форм ..	77
3.4. Определение МП и расширенного МП-автоматов. Примеры построения	88
3.5. Способы реализации синтаксических анализаторов	99
3.6. LL(k)-грамматики и пример построения LL(k)-анализатора	00
3.7. LR(k)-грамматики и пример построения LR(k)-анализатора ...	22
3.8. Грамматики предшествования	33
3.9. Пример применения алгоритм “перенос-свертка”	33
	63
Глава 4. Грамматики общего вида и машины Тьюринга	65
4.1. Грамматики общего вида и машина Тьюринга	65
4.2. Контекстно-зависимые грамматики и ленточные автоматы ..	65
4.3. Соотношение между грамматиками и языками	66
4.4. Способы реализации	66
Глава 5. Формальные методы описания перевода	67
5.1. Перевод и семантика	22
5.2. СУ-схемы	22
5.3. Транслирующие грамматики	22
5.4. Атрибутивные транслирующие грамматики	
5.5. Методика разработки описания перевода.	
5.6. Пример разработки АТ грамматики	
Глава 6. Разрработка и реализация синтаксически управляемого перевода	
6.1. L-аттрибутивные и S-аттрибутивные транслирующие грамматика	
6.2. Форма простого присваивания	
6.3. Атрибутивный перевод для LL(1) грамматик.	
6.4. S-аттрибутивный ДМП-процессор	

Приложение А. Порядок выполнения задач
Приложение В. Задания для усвоения метериала
Заключение
Библиографический список

Введение

Пособие направлено на освоение и применение теории компиляции. Приводятся примеры построения синтаксических анализаторов для регулярных и контекстно-свободных грамматик, дается реализация анализаторов на объектно-ориентированном языке программирования С# [1].

Отличительной особенностью пособия является практическое содержание, включая проектирование грамматик и реализацию распознавателей. Изложенный теоретический материал используется для проведения практических работ.

Пособие состоит из трех глав и приложения, в котором приведены основные шаги выполнения практических работ и фрагменты программных текстов. Теоретический материал и практические задания организованы по темам, в соответствии с иерархией Хомского.

В первой главе даются основные определения грамматик и языков, приводятся примеры построения синтаксических анализаторов и приемы программирования.

Основной задачей практических заданий является изучение методов проектирования грамматических правил [3-5], синтаксических анализаторов и их реализации на объектно-ориентированном языке С# [2]. В процессе выполнения практических заданий развивается умение разрабатывать грамматические правила, проводить сравнительный анализ грамматик, определять их вид в соответствии с иерархией Хомского, а также применять объектно-ориентированный подход и реализовывать грамматические правила и синтаксические анализаторы на языке программирования С#.

Во второй главе рассматривается класс клеточных автоматов, которые широко применяются в различных приложениях искусственного интеллекта и при решении дифференциальных уравнений.

В третьей главе рассмотрен класс систем, позволяющий строить различные самоподобные множества по цепочкам порожденным КС-грамматиками. например, дерево Фибоначчи. Система определяет способ генерации совокупностей бесконечных цепочек по шаблону. Все нетерминалы замещаются параллельно. Шаблон содержит терминальные символы и определяет структуру порождаемых слов.

Приводится пример фрактального анализа дерева Фибоначчи, основывающийся на определении самоподобных структур объектов.

Глава 1. Классификация грамматик и автоматов

В главе рассматриваются два способа реализации грамматик: порождающими и распознающими системами см. рис. 1.1.

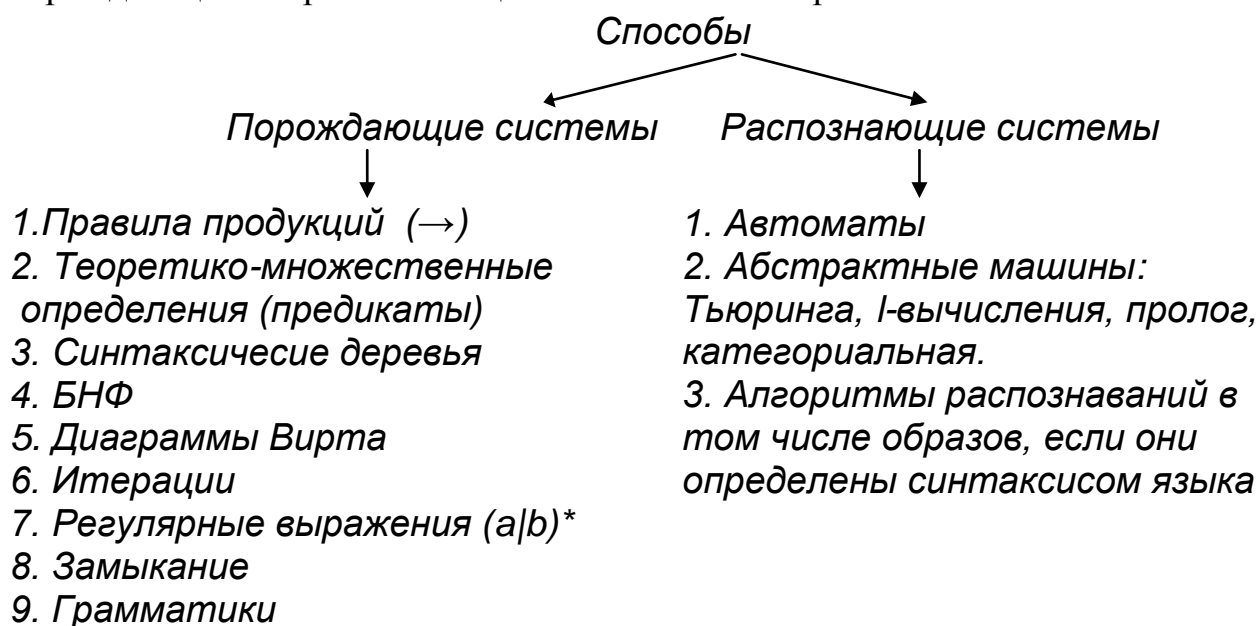


Рис. 1.1. Способы описания синтаксиса языков

Дается классификация Хомского, вводятся определения грамматик и эквивалентных им автоматов. Приводятся соотношения между типами грамматик и языков

1.1. Способы описания синтаксиса языков

Язык, будь то естественный (естественная речь) или язык программирования, задается множеством правил, образующих грамматику языка. Синтаксис языка отделяется от семантики. Появление формальных грамматик было обусловлено стремлением формализовать естественные языки и языки программирования.

Фигурные скобки используются для обозначения множеств, например $\{1,2,3\}$ обозначает множество, содержащее целые числа 1, 2 и 3.

Объединение \cup и пересечение \cap множеств определяются как обычно:

$$\{1,2,3\} \cup \{3,4,5\} = \{1,2,3,4,5\}$$

$$\{1,2,3\} \cap \{3,4,5\} = \{3\}$$

Множество A включает \supseteq множество B , если каждый элемент B является элементом A , например

$$\{3,4,5\} \supseteq \{3\}$$

Если A есть множество и мы определяем B как

$$B = \{x \mid x \in \mathbb{N}, x \text{ есть четное число}\}$$

то B будет множеством $\{2,4,8\}$ задается с помощью предиката, в данном случае это " $x \in \mathbb{N}, x$ есть четное число".

Правила описывают те последовательности слов, которые являются корректными, т.е. допустимыми предложениями языка. Грамматика позволяет осуществить грамматический анализ предложения, а значит, и явно описать его структуру.

Одно из назначений грамматики – это определить бесконечное число предложений языка с помощью конечного числа правил.

Грамматика языка может быть реализована двумя способами: в виде порождающей системы и в виде устройств, называемых распознавателями.

1. Порождающие системы. Для описания синтаксиса языков программирования наибольшее распространение получили следующие порождающие системы: формальные грамматики, форма Бекуса-Наура и ее модификации, диаграммы Вирта [3].

В порождающих системах правила грамматики задаются продукциями, состоящими из двух частей: левой и правой. В левой части, записывается определение, в правой части варианты, из которых может состоять определение. Например, предложение “он смотрит кино” может быть определено синтаксическими правилами:

предложение состоит из *подлежащее*, *группа сказуемого*
подлежащее - он
группа сказуемого - *сказуемое*, *дополнение*
сказуемое - смотрит
дополнение – кино

Синтаксис и семантика

a b c
“он смотрит кино”

В рассматриваемых правилах порядок элементов фиксирован. Для определения грамматического формализма, эквивалентного синтаксическому, введем правило замены, которое обозначается символом “→”. Тогда порождающая грамматика, для рассматриваемого предложения, примет вид:

Предложение → *подлежащее*, *группа сказуемого*
подлежащее → он
группа сказуемого → *сказуемое*, *дополнение*
сказуемое → смотрит
дополнение → кино

Выполнив правило замены “→”, получим предложение “он смотрит кино”. Предложения языка, порождаются из исходного символа применением продукций. Шаги замены можно интерпретировать как шаги построения дерева вывода см. рис. 1.2.

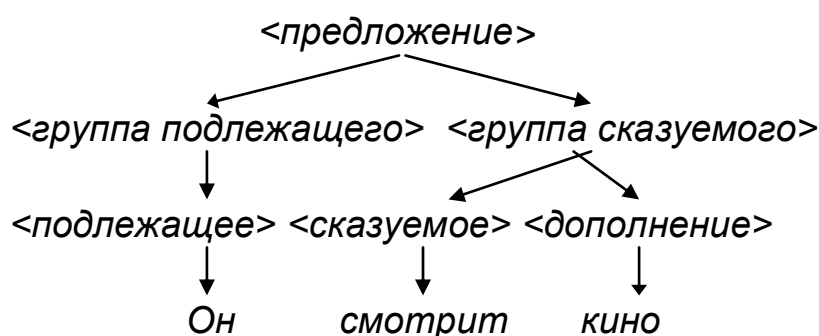


Рис. 1.2. Древовидная структура предложения

Бэкусовая нормальная форма (БНФ) или форма Бэкуса-Наура была предложена Д.Бэкусом в 1959 году и впервые применена П.Науром для описания языка Алгол-60. БНФ - *метаязык*, который используется для описания других языков.

Обозначение правил в форме записи $\xi ::= \eta$ относится к нотации БНФ. В этой нотации используются обозначения:

$::=$ - это есть, $|$ - или, $< >$ - угловые скобки, в которых записываются металингвистическая переменная, т.е. определяемое понятие;

$[]$ - факультативный элемент (необязательная часть), то есть конструкция в скобках может присутствовать или отсутствовать во фразе языка;

$\{ \}$ - множественный элемент (одно из, элемент выбора), то есть во фразе языка используется один из элементов внутри скобки.

Правило вида $\alpha ::= \beta \mid \beta\gamma$ можно представить $\alpha ::= \beta [\gamma]$.

БНФ для рассматриваемого примера.

$\langle \text{предложение} \rangle ::= \langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle$

$\langle \text{подлежащее} \rangle ::= \text{он} \mid \text{она}$

$\langle \text{группа сказуемого} \rangle ::= \langle \text{сказуемое} \rangle \langle \text{дополнение} \rangle$

$\langle \text{сказуемое} \rangle ::= \text{смотрит} \mid \text{обожает}$

$\langle \text{дополнение} \rangle ::= \text{кино}$

Используя множество правил, *выведем* или *породим* предложение по следующей схеме. Начнем с начального символа грамматики - $\langle \text{предложение} \rangle$, найдем правило, в котором $\langle \text{предложение} \rangle$ слева от $::=$, и подставим вместо $\langle \text{предложение} \rangle$ цепочку, которая расположена справа от $::=$, т.е.

$\langle \text{предложение} \rangle \Rightarrow \langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle$

Символ " \Rightarrow " означает, что один символ слева от \Rightarrow в соответствии с правилом грамматики заменяется цепочкой, находящейся справа от \Rightarrow . Обозначение \Rightarrow называют также непосредственное следование.

Заменяем синтаксическое понятие на одну из цепочек, из которых оно может состоять. Повторим процесс. Возьмем один из метасимволов в цепочке $\langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle$, например $\langle \text{подлежащее} \rangle$; найдем правило, где $\langle \text{подлежащее} \rangle$ находится слева от $::=$, и заменим $\langle \text{подлежащее} \rangle$ в исходной цепочке на соответствующую цепочку, которая находится справа от $::=$. Получим вывод:

$\langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle \Rightarrow \text{он} \langle \text{группа сказуемого} \rangle$

Полный вывод одного предложения будет таким:

$\langle \text{предложение} \rangle \Rightarrow \langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle \Rightarrow \text{он} \langle \text{группа сказуемого} \rangle \Rightarrow \text{он} \langle \text{сказуемое} \rangle \langle \text{дополнение} \rangle \Rightarrow \text{он} \text{смотрит} \langle \text{дополнение} \rangle \Rightarrow \text{он} \text{смотрит} \text{кино}$

Этот вывод предложения запишем сокращенно, используя символ \Rightarrow^+ , который означает, что цепочка выводима *нетривиальным способом*:

$\langle \text{предложение} \rangle \Rightarrow^+ \text{он} \text{смотрит} \text{кино}$

Две последовательности связаны отношением \Rightarrow^* , когда вторая получается из первой применением некоторой последовательности продукции $\alpha \Rightarrow^* \alpha_m$ тогда, и только тогда, когда $\alpha \Rightarrow^+ \alpha_m$.

На каждом шаге можно заменить любую метапеременную. В приведенном выше выводе всегда заменялся самый левый из них. Если в процессе вывода цепочки правила грамматики применяются только к самому левому нетерминалу, говорят, что получен *левый вывод* цепочки. Аналогично определяется правый вывод.

Рассматриваемая грамматика определяет несколько предложений (цепочек) языка:

он смотрит кино он обожает кино
она смотрит кино она обожает кино

Диаграммы Вирта позволяют визуально представить правила порождения. Каждому правилу соответствует диаграмма.

Элементы ИЛИ обозначаются разветвлениями и вершинами, элементы И последовательными вершинами графа, метасимволы $\{...\}$ циклом. Вершины могут быть двух типов: терминальные – кружки, и нетерминальные – прямоугольники. Каждому правилу соответствует диаграмма см. рис. 1.3.

Правила- И, цикл, ИЛИ: $T \rightarrow \{ F \} | (E)$

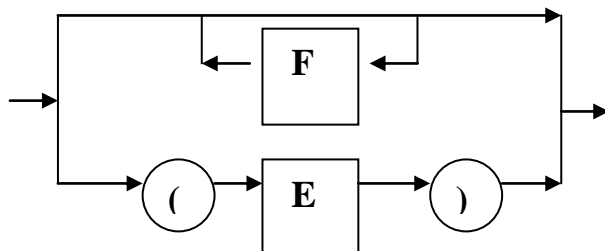


Рис. 1.3. Диаграммы Вирта

Итерационная форма описания. Вводится операция итерации – обозначается парой круглых скобок со звездочкой.

Итерация вида $(a)^*$ определяется как множество, включающее цепочки всевозможной длины, построенные с использованием символа **a**.

$$(a)^* = \{a, aa, aaa, aaaa, \dots\}$$

Определим понятия, используемые при определении грамматик формальных языков и распознавателей.

Алфавит (или словарный состав) – конечное множество символов. Например, *алфавит* $V_1 = \{a,b\}$ – содержит два символа, *алфавит* $V_2 = \{01,11,10,00\}$ – содержит четыре символа.

Цепочкой символов α в алфавите V называется любая конечная последовательность символов этого алфавита.

Пустой цепочкой символов ϵ называется цепочка, не содержащая ни одного символа.

Длина цепочки α – число составляющих ее символов, обозначается $|\alpha|$. Например, если цепочка символов $\alpha = \text{defgabck}$, то длина α равна 8, $|\alpha| = 8$. Длина ϵ равна 0, $|\epsilon| = 0$.

Замыкание К्लीни (или звезда К्लीни) в математической логике и информатике – унарная операция над множеством строк либо символов. Замыкание К्लीни множества V обозначается V^* . Широко применяется в регулярных выражениях, на примере которых было введено Стивеном К्लीни для описания некоторых автоматов.

Множество всех строк (включая пустую), которые могут быть построены из **символов** алфавита V , называется **замыканием** V , и обозначается V^* .

Пусть V – множество строк, тогда V^* – минимальное надмножество множества V , которое содержит ε (пустую строку) и замкнуто относительно конкатенации.

Надмножество (superset) – множество, подмножеством которого является данное. Иначе супермножество, объемлющее множество (множество множеств).

V^* (замыкание V) – обозначается множество всех цепочек составленных из символов алфавита V , включая пустую цепочку ε . Это также множество всех строк, полученных конкатенацией ε или более строк из V . Например, если $V = \{0, 1\}$, то $V^* = \{\varepsilon, 0, 1, 00, 11, 01, 10, 000, 001, 011, \dots\}$.

V^+ – обозначается множество всех цепочек составленных из алфавита V , исключая пустую цепочку ε . Следовательно, $V^* = V^+ \cup \{\varepsilon\}$.

Декартовым произведением $A \times B$ множеств A и B называется множество $\{(a, b) \mid a \in A, b \in B\}$.

\emptyset – обозначается пустое множество.

2. Распознающие системы (автоматы). Грамматика может быть реализована в виде некоторого алгоритма, производящего действия и отвечающего на вопрос, принадлежит ли данная фраза языку.

Функцией перехода δ называется отношение, управляющее действиями распознавателя. Функция перехода δ определяет для каждой пары (входной символ, состояние) множество состояний, в которых он будет находится на следующем шаге.

Например, принадлежит ли предложение “он смотрит кино” языку заданному рассмотренными грамматическими правилами? Для этого необходимо построить распознаватель.

1.2. Классификация Хомского

Синтаксис языка L можно специфицировать, пользуясь системой изображения множеств, например:

$$L = \{0^n 1^n \mid n \geq 0\}$$

Язык L включает строки, состоящие из нуля или нулей и того же числа последующих единиц. Пустая строка включается в язык.

Такое задание синтаксиса намного проще синтаксисов большинства языков специфицируемых с помощью *грамматики*.

Определение 1. Грамматика $G = (T, V, P, S_0)$,

где T – конечное множество терминальных символов (терминалов) алфавита;

V – конечное множество нетерминальных символов алфавита, не пересекающихся с T , $T \cap V = \emptyset$,

S_0 – начальный символ (или аксиома), $S_0 \in V$;

P – конечное множество правил порождения (продукций), $P = (T \cup V)^+ \times (T \cup V)^*$. Элемент (α, β) множества P называется *правилом вывода* и записывается в виде $\alpha \rightarrow \beta$.

Цепочка $\beta \in (T \cup V)^*$ непосредственно выводима ($\alpha \rightarrow \beta$) из цепочки $\alpha \in (T \cup V)^+$ в грамматике G , если $\alpha = \xi_1 \gamma \xi_2$, $\beta = \xi_1 \delta \xi_2$, где $\xi_1, \xi_2, \delta \in (T \cup V)^*$, $\gamma \in (T \cup V)^+$ и правило вывода $\gamma \rightarrow \delta$ содержится в P .

Цепочка $\beta \in (T \cup V)^*$ выводима из цепочки $\alpha \in (T \cup V)^+$ в грамматике G , (обозначим $\alpha \Rightarrow \beta$), если существуют цепочки $\gamma_0, \gamma_1, \dots, \gamma_n$ ($n \geq 0$), такие, что $\alpha = \gamma_0 \rightarrow \gamma_1 \rightarrow \dots \rightarrow \gamma_n = \beta$. Последовательность $\gamma_0, \gamma_1, \dots, \gamma_n$ называется *выводом длины n* .

Правила вывода с одинаковыми левыми частями $\alpha \rightarrow \beta_1, \alpha \rightarrow \beta_2, \dots, \alpha \rightarrow \beta_n$, записывают сокращенно $\alpha \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$, где $\beta_i, i = 1, 2, \dots, n$ и называют *альтернативой* правил вывода из цепочки α .

Для обозначения терминалов грамматики мы будем употреблять прописные буквы (или строки прописных букв), а для обозначения не терминалов — заглавные буквы (или строки заглавных букв).

Примеры грамматик.

Пример 1. Для примера “он смотрит кино” имеем грамматику $G = (\{\text{“он”}, \text{“кино”}, \text{“смотрит”}\}, \{\text{Предложение, подлежащее, группа сказуемого, сказуемое, дополнение}\}, P, S_0)$, где $S_0 = \text{Предложение}$, P состоит из правил $P = \{p_1, p_2, p_3, p_4, p_5\}$;

p_1 : Предложение \rightarrow подлежащее, группа сказуемого

p_2 : подлежащее \rightarrow он

p_3 : группа сказуемого \rightarrow сказуемое, дополнение

p_4 : сказуемое \rightarrow смотрит

p_5 : дополнение \rightarrow кино

Пример 2. $G = (\{c, d\}, \{B, S_0\}, P, S_0)$, где P состоит из правил

$S_0 \rightarrow cBd$

$cB \rightarrow ccBd$

$B \rightarrow \varepsilon$

Цепочка $ccBdd$ непосредственно выводима из cBd в грамматике G .

Цепочка $cccBddd$ в грамматике G выводима, т.к. существует вывод $S_0 \Rightarrow cBd \Rightarrow ccBdd \Rightarrow cccBddd \Rightarrow cccddd$. Длина вывода равна 4.

Определение 2. Языком, порождаемым грамматикой $G = (T, V, P, S_0)$, называется множество $L(G) = \{\alpha \in T^* \mid S_0 \Rightarrow^* \alpha\}$.

Другими словами, $L(G)$ - это все цепочки в алфавите T , которые выводимы из S_0 с помощью P .

Цепочка $\alpha \in (T \cup V)^*$, для которой $S_0 \Rightarrow^* \alpha$ (то есть цепочка, которая может быть выведена из начального символа), называется *сентенциальной формой* в грамматике G . Язык, порождаемый грамматикой, можно определить как множество терминальных сентенциальных форм.

Например, пусть задана грамматика $G_1 = (\{0, 1\}, \{A, S_0\}, P_1, S_0)$, состоящая из правил P_1 : $S_0 \rightarrow 0A1$

$0A \rightarrow 00A1$

$A \rightarrow \varepsilon$

Тогда язык, порождаемый грамматикой G_1 , это язык $L(G_1) = \{0^n 1^n \mid n > 0\}$.

Грамматик G_1 и G_2 называются *эквивалентными*, если $L(G_1) = L(G_2)$. Например, пусть задана грамматика $G_2 = (\{0, 1\}, \{S\}, P_2, S)$, где

P_2 : $S \rightarrow 0S1 \mid 01$, тогда

грамматики G_1 и G_2 эквивалентны, т.к. обе порождают язык $L(G_1) = L(G_2) = \{0^n 1^n \mid n > 0\}$.

Грамматика G_2 и G_3 называются *почти эквивалентными* $L(G_2) = L(G_3)$, если $L(G_2) = L(G_3) \cup \{\varepsilon\}$. Другими словами, грамматики почти эквивалентны, если языки, ими порождаемые, отличаются не более чем на ε .

Определим грамматику $G_3 = (\{0,1\}, \{S\}, P_3, S)$ с правилами

$$P_3: S \rightarrow 0S1 \mid \varepsilon$$

Тогда грамматики G_2 и G_3 *почти эквивалентны*, так как $L(G_2) = \{0^n 1^n \mid n > 0\}$, а $L(G_3) = \{0^n 1^n \mid n \geq 0\}$, т.е. $L(G_3)$ состоит из всех цепочек языка $L(G_2)$ и пустой цепочки, которая в $L(G_2)$ не входит.

Классификация Хомского в таб.1.1. дает представление о 4 типах грамматик, располагающихся по иерархии от 0 до 3 в порядке убывания их общности, и о эквивалентных формальных автоматах.

Таб. 1.1. Классификация Хомского

Тип	Ограничения на правила P	Пример правил грамматики $G = (T, V, P, S_0)$	Пример языка	Автомат или абстрактная машина
0	Общего вида, не накладывається никаких ограничений	$S_0 \rightarrow CD, C \rightarrow 0CA, C \rightarrow 1CB, AD \rightarrow 0D, BD \rightarrow 1D, A0 \rightarrow 0A, A1 \rightarrow 1A, B0 \rightarrow 0B, B1 \rightarrow 1B, C \rightarrow \varepsilon, D \rightarrow \varepsilon$	$\{\omega\omega \omega \in \{0,1\}^*, \text{ язык состоит из цепочек четной длины, из 0 и 1}\}$	Машина Тьюринга
1	Контекстно-зависимая (КЗ) $\alpha \rightarrow \beta, \alpha \leq \beta $	$S_0 \rightarrow 0B0$ $B \rightarrow 1 \mid 0BC$ $C0 \rightarrow 100$ $C1 \rightarrow 1C$	$\{0^n 1^n 0^n n \geq 1\}$	Машина Тьюринга с конечной лентой
2	Контекстно-свободная (КС) $A \rightarrow \alpha$	$S_0 \rightarrow AS_0B \mid AB$ $A \rightarrow 0$ $B \rightarrow 1$ LL(k) грамматика LR(k) грамматика	$\{0^n 1^n n \geq 1\}$	с магазинной памятью (МП-автомат: LL-разбор, РМП-автомат: LR-разбор) Распознаватели: LL(1) LR(1) подкласс LR(1) грамматики предшествования
3	Регулярная $A \rightarrow bV$ $B \rightarrow a$	$S_0 \rightarrow 0 \mid 0A$ $A \rightarrow 0B$ $B \rightarrow 1 \mid 1A$	$\{0(01)^n n \geq 0\}$	Конечный автомат (КА)

Грамматики классифицируются по виду правил вывода P. В правилах вывода P большими латинскими буквами обозначают нетерминальные символы, маленькими латинскими буквами - терминальные.

Автомат эквивалентен грамматике, если он воспринимает весь порожденный грамматикой язык и только этот язык.

Классификация Хомского дает представление об общности грамматик, и об эквивалентных им формальных автоматах.

Глава 2. Автоматные грамматики и конечные автоматы

2.1. Определение автоматных грамматик и конечных автоматов

Определение 3. Грамматика типа 3 или автоматная (регулярная) A -грамматика – это грамматика $G = (T, V, P, S_0)$, у которой правила порождения P имеют вид по классификации Хомского $A \rightarrow bV$ (праволинейное правило) или $V \rightarrow a$ (заключительное правило), где $A, V \in V$, $a, b \in T$. Каждое правило такой грамматики содержит единственный нетерминал в левой части, всегда один терминал в правой части, за которым может следовать один нетерминал. Такую грамматику также называют *праволинейной*.

Грамматика $G = (T, V, P, S_0)$ называется *леволинейной*, если каждое правило из P имеет вид по классификации Хомского $A \rightarrow Vb$ либо $V \rightarrow a$, где $A, V \in V$, $a, b \in T$.

Таким образом, грамматику типа 3 можно определить как праволинейную либо как леволинейную.

Выбор определения не влияет на множество языков, порождаемых грамматиками этого класса, поскольку доказано, что множество языков, порождаемых праволинейными грамматиками, совпадает с множеством языков, порождаемых леволинейными грамматиками.

Автоматной грамматике эквивалентен простейший распознаватель – недетерминированный конечный автомат. Автомат представляет собой устройство, у которого отсутствует вспомогательная память.

Определение 4. Конечный автомат (КА) – это пятерка объектов

$KA = (Q, \Sigma, \delta, q_0, F)$, где

Q - конечное множество состояний;

Σ - конечный алфавит входных символов;

δ - функция переходов, задаваемая отображением

$$\delta: Q \times \Sigma \rightarrow Q,$$

где Q - конечное множество подмножеств множества Q ;

q_0 - начальное состояние автомата, $q_0 \in Q$;

$F \subseteq Q$ - множество заключительных состояний.

В каждый момент времени КА находится в некотором состоянии $q \in Q$ и читает поэлементно последовательность символов $a_k \in \Sigma$, записанную на конечной ленте. При этом либо читающая головка машины движется в одном направлении (слева направо), либо лента перемещается (справа налево) рис.2.1. при неподвижной читающей головке.



Рис. 2.1. Схема конечного автомата
 Входная лента: Входные символы $a_i \in \Sigma$

Если автомат в состоянии q_i читает символ a_k и определена функция перехода $\delta(q_i, a_k) = q_j$, то автомат воспринимает символ a_k и переходит в состояние q_j для обработки следующего символа.

Определение 5. Конфигурация КА – это пара множества $(q, \omega) \in Q \times \Sigma^*$, где $q \in Q$, $\omega \in \Sigma^*$. Конфигурация (q_0, ω) называется *начальной*, а (q, ε) , где $q \in F$, – *заключительной*.

Определим бинарное отношение \vdash на конфигурациях, соответствующее одному такту работы КА. Если $q' \in \delta(q, a)$, то $(q, a\omega) \vdash (q', \omega)$ для всех $\omega \in \Sigma^*$.

Пусть $\{C\}$ – множество конфигураций.

1. $C \vdash^0 C'$ означает, что $C = C'$
2. $C_0 \vdash^k C_k$, если существует последовательность конфигураций C_1, C_2, \dots, C_{k-1} , для $k \geq 1$, в которых $C_i \vdash C_{i+1}$ для $0 \leq i < k$.
3. $C \vdash^+ C'$ означает, что $C \vdash^k C_k$ для некоторого $k \geq 1$, а $C \vdash C'$ означает, что $C \vdash^k C'$ для $k \geq 1$.

Конечный автомат $KA = (Q, \Sigma, \delta, q_0, F)$ распознает входную цепочку $\omega \in \Sigma^*$, если $(q_0, \omega) \vdash^k (q, \varepsilon)$ для $q \in F$.

Языком $L(KA)$, распознаваемым КА называется множество входных цепочек,

$$L(KA) = \{\omega \in \Sigma^* \mid (q_0, \omega) \vdash^k (q, \varepsilon), \text{ где } q \in F\}$$

КА называется *недетерминированным*, если для каждой конфигурации существует конечное множество всевозможных следующих шагов, любой из которых КА может сделать, исходя из этой конфигурации.

КА называется *детерминированным*, если для каждой конфигурации существует не более одного следующего шага.

Для любого конечного недетерминированного автомата можно построить ему эквивалентный детерминированный автомат.

Класс языков, распознаваемый конечными автоматными, совпадает с классом языков, порождаемых автоматными грамматиками и наоборот.

Утверждение 1. Пусть задана автоматная грамматика $G = (T, V, P, S_0)$, тогда существует такой (недетерминированный) конечный автомат $KA = (Q, \Sigma, \delta, q_0, F)$, что $L(KA) = L(G)$. КА строится следующим образом:

1. Входные алфавиты конечного автомата Σ и автоматной грамматики T совпадают, $\Sigma = T$.
2. $Q = V \cup \{q_f\}$, где q_f – заключительное состояние конечного автомата.

3. $q_0 = S_0$.
4. Если $S_0 \rightarrow \varepsilon \in P$, то $F = \{ S_0, q_f \}$, в противном случае $F = \{ q_f \}$.
5. Функция переходов КА определяется следующим образом:
 - $q_f \in \delta(B, a)$, если $B \rightarrow a \in P$, $B \in V$, $a \in \Sigma$;
 - если $B \rightarrow aC \in P$, то $C \in \delta(B, a)$;
 - $\delta(q_f, a) = \emptyset$ для всех $a \in \Sigma$.

Пример. Пусть задан регулярный язык $L = \{0(10)^n \mid n \geq 0\}$. Построить автоматную грамматику $G = (T, V, P, S_0)$ для заданного языка L и привести пример вывода строки. Используя грамматику G , построить

КА $= (\Sigma, Q, \delta, q_0, F)$ и привести пример конфигурации КА.

1. Построение грамматики. $L = \{0(10)^n \mid n \geq 0\}$, то 0, 010, 01010 и т.д. - этот язык порожден регулярной грамматикой $G = (T, V, P, S_0)$, где $T = \{0, 1\}$, $V = \{S_0, A, B\}$, $P = \{S_0 \rightarrow 0, S_0 \rightarrow 0A, A \rightarrow 1B, B \rightarrow 0, B \rightarrow 0A\}$. Пример вывода цепочки $S_0 \Rightarrow 0A \Rightarrow 01B \Rightarrow 010A \Rightarrow 0101B \Rightarrow 01010$.

2. Построение КА. Воспользуемся утверждением 1, тогда

КА $= (\{0, 1\}, \{S_0, A, B, q_f\}, \delta, S_0, q_f)$

$\delta(S_0, 0) = \{A, q_f\}$

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

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

Пример конфигурации КА: $S_0 01010 \mid A 1010 \mid B 010 \mid A 10 \mid B 0 \mid q_f, \varepsilon$

2.2. Способы задания конечных автоматов

На рис. 2.2. приведена диаграмма переходов КА для рассматриваемого примера.

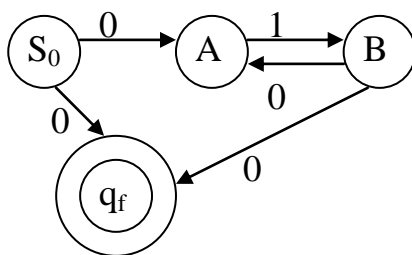


Рис. 2.2. Диаграмма переходов КА

Конечный автомат можно задать в виде таблицы переходов и диаграммы (графа) переходов.

В таблице переходов двум аргументам ставится в соответствие одно значение. В таб. 2.1 приведена таблица переходов КА для рассматриваемого примера.

Таб. 2.1. Таблица переходов для КА

	0	1
S_0	$\{A, q_f\}$	
A		$\{B\}$
B	$\{A, q_f\}$	
q_f		

Диаграммой переходов КА называется неупорядоченный граф, удовлетворяющий условиям:

- каждому состоянию q соответствует некоторая вершина, отмеченная его именем;
- диаграмма переходов содержит дугу из состояния q_k в состояние q_n , отмеченную символом a , если $q_k \in \delta(q_n, a)$. Дуга может быть помечена множеством символов, переводящих автомат из состояния q_k в состояние q_n ;
- вершины, соответствующие заключительным состояниям $q_f \in F$, отмечаются двойным кружком.

2.3. Свойства регулярных языков: лемма о накачке

Теорема, называемая “лемма о накачке” утверждает, что все достаточно длинные слова регулярного языка можно *накачать*, то есть **повторить** внутреннюю часть слова сколько угодно раз, производя новое слово, также принадлежащее языку. Лемма описывает существенное свойство всех регулярных языков и служит инструментом для доказательства нерегулярности некоторых языков. Она является одной из нескольких лемм о накачке.

Два способа:

1. построить распознаватель для заданного языка - конечный автомат;
2. выделить цепочку y^i для всех $i \geq 0$, $xy^iz \in L$, на x и z ограничений не накладывается (то есть, нет цепочки символов y^i , то не регулярный язык).

Свойство замкнутости регулярных языков, позволяет минимизировать построение автомата:

1. строить распознаватели для одних языков, построенных из других с помощью операций;
2. определить, что два различных автомата определяют один язык.

Пример 1. Пусть КА (см. рис.) распознает строку **abcbcd**. Поскольку длина её превышает число состояний, существуют повторяющиеся состояния: q_1 и q_2 .

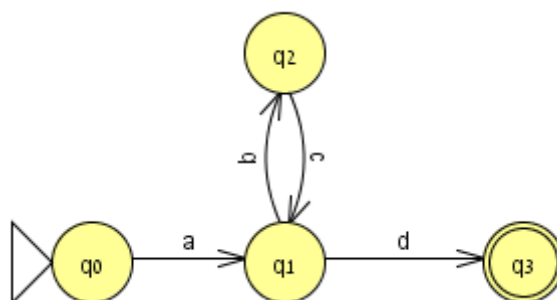


Рис. 2.3. Диаграмма переходов КА

Поскольку подстрока **bcbc** строки **abcbcd** проводит автомат по переходам из состояния q_1 обратно в q_1 , эту строку можно повторять сколько угодно раз, и КА всё равно будет её принимать, например строки **abcbcbcbcd** и **ad**.

В терминах леммы о накачке, строка **abcbcd** разбивается на часть $x = a$, часть $y = bc$ и часть $z = bcd$.

Заметим, что можно разбить её различными способами, например $x = \mathbf{a}$, $y = \mathbf{bcbs}$, $z = \mathbf{d}$, и все условия будут выполнены, кроме $|xy| \leq p$, где p – длина накачки.

Длина накачки p принимается более длины самого длинного слова языка.

Неформальное утверждение: слово w языка L длины по меньшей мере p , (где p константа, называемая длиной накачки, зависит лишь от L) можно разделить на три подцепочки, $w = xyz$, так что среднюю часть, y (непустую), можно повторить произвольное число раз (включая ноль, **то есть удалить**) и получить строку из L . Этот процесс повторения называется «накачкой».

Причем, длина xy не превысит p , ограничивая способы деления строки w .

Конечные языки удовлетворяют требованиям леммы о накачке тривиально определяя m длиной максимальной строки из языка плюс один.

Пример 2. Рассмотрим язык $L_{01} = \{0^n 1^n \mid n \geq 1\}$, состоящий из всех цепочек вида 01, 0011, 000111 и так далее, содержащий один или несколько нулей, за которыми следует такое же количество единиц.

Утверждается, что язык L_{01} нерегулярен.

Если бы L_{01} был регулярным языком, то допускался некоторым ДКА, имеющим какое-то число состояний k . Пусть на вход ДКА поступает k нулей. Он находится в некотором состоянии после чтения каждого из $k+1$ префиксов входной цепочки, т.е. $\epsilon, 0, 00, \dots, 0^k$. Поскольку есть только k различных состояний, например, $0^i, 0^j$, автомат должен находиться в одном и том же состоянии.

Прочитав i или j нулей ДКА получает на вход 1. По прочтении i единиц он должен допустить вход, если ранее получил i нулей, и отвергнуть его, получив j нулей. Но в момент поступления 1, автомат не способен вспомнить какое число нулей i, j было принято. Следовательно, он может работать неправильно.

Формальное утверждение.

Пусть L регулярный язык. Тогда существует целое $p \geq 1$ зависящее только от L , такое что цепочка w из L длины по меньшей мере p может быть записана как $w = xyz$, где y – это подцепочка, которую можно накачать (удалить или повторить произвольное число раз, так что результат останется в L), при этом:

1. $|y| \geq 1$, цикл y должен быть накачан хотя бы длиной 1;
2. $|xy| \leq p$, цикл должен быть в пределах первых p символов;
3. для всех $i \geq 0$, $xy^i z \in L$, на x и z ограничений не накладывается.

$$\forall L \subseteq \Sigma^* (\text{regular}(L) \Rightarrow (\exists p \geq 1 (\forall w \in L (|w| \geq p) \Rightarrow (\exists x, y, z \in \Sigma^* (w = xyz \Rightarrow (|y| \geq 1 \wedge |xy| \leq p \wedge \forall i \geq 0 (xy^i z \in L)))))))$$

Пример 3. Нерегулярность языка $L = \{a^n b^n \mid n \geq 0\}$ над алфавитом $\Sigma = \{a, b\}$ можно показать следующим образом.

Пусть w, x, y, z, p , и i заданы соответственно формулировке леммы выше. Пусть w из L задаётся как $w = a^p b^p$. По лемме о накачке, существует разбиение $w = xyz$, где $|xy| \leq p$, $|y| \geq 1$, такое что $xy^i z$ принадлежит L для любого $i \geq 0$. Если

допустить, что $|xy|=p$, а $|z|=p$, то xu — это первая часть w , состоящая из p последовательных экземпляров символа a .

Поскольку $|y| \geq 1$, она содержит по меньшей мере одну букву a , а xu^2z содержит больше букв a чем b . Следовательно, xu^2z не в языке L (заметим, что любое значение $i \neq 1$ даст противоречие). Достигнуто противоречие, поскольку в этом случае накачанное слово не принадлежит языку L . Предположение о регулярности L неверно и L — не регулярный язык.

Доказательство леммы о накачке

Для каждого регулярного языка существует конечный автомат (КА), распознающий этот язык.

Если язык конечен, то результат можно получить немедленно, задав длину накачки p более длины самого длинного слова языка: в этом случае нет ни одной цепочки в языке длиннее p , и утверждение леммы выполняется безусловно.

Если регулярный язык бесконечен, то существует минимальный КА, распознающий его. Число состояний этого КА и принимается за длину накачки p . Если длина цепочки превышает p , то хотя бы одно состояние при обработке цепочки повторяется (назовём его S). Переходы из состояния S и обратно соответствуют некоторой цепочки. Эту цепочку обозначим u по условиям леммы, и поскольку автомат примет строку как без части u , так и с повторяющейся частью u , условия леммы выполнены.

Доказательство. Пусть L — регулярный язык, тогда $L = L(\text{ДКА})$ для некоторого ДКА. Пусть ДКА имеет n состояний. Рассмотрим произвольную цепочку w длиной не менее n . и скажем, $w = a_1a_2a_3\dots a_m$, где $m \geq n$ и каждый a_i есть входной символ. Для $i = 0, 1, 2, \dots, n$ определим состояние p_i как $\delta(q_0, a_1a_2a_3\dots a_i)$, где δ — функция переходов автомата, q_0 — его начальное состояние.

Заметим, что $p_0 = q_0$.

Рассмотрим $n+1$ состояний p_i при $i = 0, 1, 2, \dots, n$. Поскольку автомат имеет n различных состояний, то найдутся два разных целых числа i и j ($0 \leq i < j \leq n$), при которых $p_i = p_j$. Теперь разобьём цепочку w на xuz .

1. $x = a_1a_2\dots a_i$
2. $u = a_{i+1}a_{i+2}\dots a_j$
3. $z = a_{j+1}a_{j+2}\dots a_m$

Таким образом, x приводит в состояние p_i , u — из p_i обратно в p_i (так как $p_i = p_j$), а z — это остаток цепочки w . Взаимосвязи между цепочками и состояниями показаны на рис. Заметим, что цепочка x может быть пустой при $i=0$, а z — при $j=n=m$. Однако цепочка u не может быть пустой, поскольку i строго меньше j .

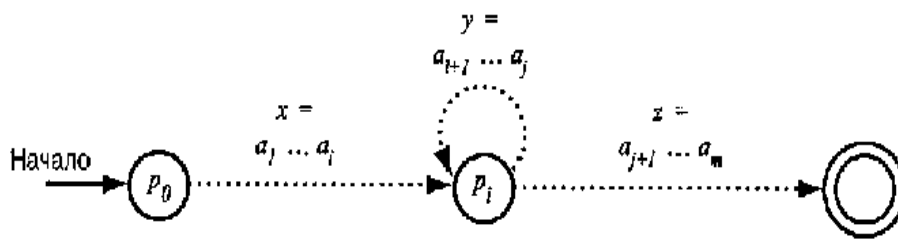


рис. 4.1. Каждая цепочка, длина которой больше числа состояний автомата, приводит к повторению некоторого состояния

Рис. 2.4. Каждая цепочка, длина которой больше числа состояний автомата приводит к повторению некоторого

На вход автомата поступает цепочка xu^kz для любого $k \geq 0$.

1. При $k = 0$ автомат переходит из начального состояния q_0 (которое есть также p_0) в p_i прочитав x . Поскольку $p_i = p_j$, то z переводит автомат из p_i в допускающее состояние (рис. 4.1.).

2. Если $k > 0$, то по x автомат переходит из q_0 в p_i , затем читая u^k , k раз циклически проходит через p_j , и, наконец, по z переходит в допускающие состояния.

3. Для любого $k \geq 0$ цепочка xu^kz также допускается автоматом, то есть принадлежит языку L .

2.4. Способы преобразования НКА в ДКА

1. Построение НКА из автоматов.

???????????????????????????????? см. лаб.

2. Детерминированный конечный автомат ДКА $= (Q, \Sigma, \delta, q_0, F)$, если множество $\delta(q, a)$ содержит не более одного состояния для любых $q \in Q, a \in \Sigma$. Рассмотрим алгоритм построения ДКА по недетерминированному КА (НКА).

Основная идея построения ДКА по НКА заключается в том, что после обработки отдельной входной цепочки состояние ДКА будет представлять собой множество всех состояний НКА, которые он может достичь из начальных состояний после применения данной цепочки.

Переходы ДКА можно получить из НКА, вычисляя множества состояний, которые могут следовать после данного множества при различных входных символах.

Допустимость цепочки определяется исходя из того, является ли последнее детерминированное состояние, которого достиг ДКА, множеством недетерминированных состояний, включающим хотя бы одно допускающее состояние.

Вход: НКА $= (Q, \Sigma, \delta, q_0, F)$

Выход: ДКА $= (Q', \Sigma, \delta', q_0', F')$, у которого $L(\text{НКА}) = L(\text{ДКА})$.

1. Входные алфавиты НКА и ДКА совпадают.

2. Множество Q' содержит множество всех подмножеств множества Q (булеан множества Q). Если $|Q| = n$, то $|Q'| = 2^n$.

3. Для каждого множества $S \subseteq Q$ и каждого входного символа $a \in \Sigma$, выполняется $\delta'(S, a) = \bigcup_{p \in S} \delta(p, a)$, т.е. чтобы найти $\delta'(S, a)$ необходимо для каждого состояния p из S найти состояния, в которые можно попасть из p для входного $a \in \Sigma$. Затем взять объединение множеств всех состояний по всем состояниям p .
4. F' есть множество подмножеств S множества Q , $S \cap F \neq \emptyset$, то есть F' состоит из всех множеств состояний Q НКА, содержащего хотя бы одно заключительное состояние.
5. $q_0' = q_0$.
6. Не все состояния из Q' достижимы из начального состояния (состояние q достижимо из начального состояния q_0 , если имеется путь из состояния q_0 в состояние q). Недостижимые состояния надо исключить из множества Q' .
7. Определить ДКА.

Построим ДКА, допускающий язык, определяемый НКА для рассматриваемого примера рис. 2.2.

1. Пусть входные алфавиты НКА и ДКА совпадают.
2. Булеан множества без пустого подмножества $2^4 - 1 = 15$, $Q' = \{\{S\}, \{A\}, \{B\}, \{q_f\}, \{S,A\}, \{S,B\}, \{S,q_f\}, \{A,B\}, \{A,q_f\}, \{B,q_f\}, \{S,A,B\}, \{S,A,q_f\}, \{A,B,q_f\}, \{S,B,q_f\}, \{S,A,B,q_f\}\}$

3. Найти $\delta'(S, a)$:

$$\begin{aligned}
 \delta'(\{B\}, 0) &= \{A, q_f\} & \delta'(\{S\}, 0) &= \{A, q_f\} & \delta'(\{A\}, 1) &= \{B\} \\
 \delta'(\{S,B\}, 0) &= \{A, q_f\} & \delta'(\{S,A\}, 0) &= \{A, q_f\} & \delta'(\{S,A\}, 1) &= \{B\} \\
 \delta'(\{S,q_f\}, 0) &= \{A, q_f\} & \delta'(\{B,A\}, 0) &= \{A, q_f\} & \delta'(\{B,A\}, 1) &= \{B\} \\
 \delta'(\{B,q_f\}, 0) &= \{A, q_f\} & \delta'(\{A,q_f\}, 1) &= \{B\} & \delta'(\{S,B,A\}, 1) &= \{B\} \\
 \delta'(\{S,B,A\}, 0) &= \{A, q_f\} & \delta'(\{S,B,q_f\}, 0) &= \{A, q_f\} & \delta'(\{B,A,q_f\}, 0) &= \{A, q_f\} \\
 \delta'(\{B,A,q_f\}, 1) &= \{B\} & \delta'(\{S,A,q_f\}, 0) &= \{A, q_f\} & \delta'(\{S,A,q_f\}, 1) &= \{B\} \\
 \delta'(\{S,B,A,q_f\}, 0) &= \{A, q_f\} & \delta'(\{S,B,A,q_f\}, 1) &= \{B\}
 \end{aligned}$$

4. $F' = \{\{q_f\}, \{A,q_f\}, \{S,q_f\}, \{B,q_f\}, \{S,B,q_f\}, \{B,A,q_f\}, \{S,A,q_f\}, \{S,B,A,q_f\}\}$
Каждое из подмножеств F' является заключительным состоянием. При определении достижимого заключительного состояния выбирается соответствующее подмножество.

5. $S = S_0$

6. Достижимыми состояниями в ДКА являются $\{S\}$, $\{B\}$ и $\{A,q_f\}$, остальные состояния удаляются.

Таб. 2.2. Таблица переходов для ДКА

	0	1
$p_0 = \{S\}$	p_1	
$p_1 = \{A, q_f\}$		p_2
$p_2 = \{B\}$	p_1	

7. ДКА = ($\{\{S\}, \{B\}, \{A, q_f\}\}$, $\{0, 1\}$, δ' , S , $\{\{A, q_f\}\}$), где

$$\delta'(\{B\}, 0) = \{A, q_f\} \quad \delta'(\{S\}, 0) = \{A, q_f\}$$

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

ДКА представляется следующим графом:

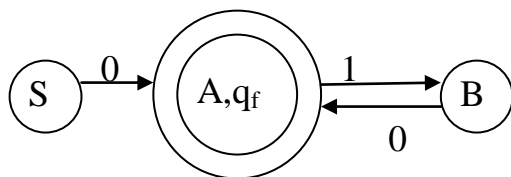


Рис. 2.5. Диаграмма переходов ДКА

Приведенный алгоритм можно использовать и для перехода от недетерминированной А-грамматики к детерминированной.

2.5. Способы реализации: по заданному регулярному выражению составные автоматы

Реализовать на языке C# КА см. рис. 2.5. Пользователь вводит цепочку символов в поле редактирования и получает ответ: “Да” – цепочка символов принадлежит языку, автомат распознал цепочку или “Нет” – цепочка символов не принадлежит заданному языку L.

Пример проектирования грамматики и распознавателя

1. Постановка задачи.

Задан язык $L = \{c\omega n \mid \omega \in \{+, d, -, k, f\}^*\}$, где ω обозначает множество всех цепочек, составленных из символов $\{+, d, -, k, f\}^*$. Цепочки начинаются с символа c, заканчиваются символом n.

Спроектировать грамматику $G = (T, V, P, S_0)$ для заданного языка L, привести пример вывода цепочки.

Используя грамматику G, построить КА $= (Q, \Sigma, \delta, q_0, F)$ и привести пример конфигурации КА. Построить диаграмму переходов КА. Определить свойства КА, если это НКА реализовать алгоритм преобразования НДКА в ДКА.

1.a. Спроектировать грамматику $G = (T, V, P, S_0)$, где

$T = \{c, d, k, f, n, +, -\}$ множество терминальных символов

$V = \{S_0, Q\}$ множество нетерминальных символов

$P = \{S_0 \rightarrow Qn, Q \rightarrow c \mid Q+ \mid Q- \mid Qd \mid Qk \mid Qf\}$

1.b. Определить свойства грамматики. Грамматика является левوليнейной, бесконечной. Пример вывода цепочки $S_0 \Rightarrow Q-n \Rightarrow Qf-n \Rightarrow Qkf-n \Rightarrow Qdkf-n \Rightarrow Q-dkf-n \Rightarrow Qd-dkf-n \Rightarrow Q+d-dkf-n \Rightarrow c+d-dkf-n$.

1.c. Используя грамматику G построить КА.

$KA = (\{S_0, Q, q_f\}, \{c, d, k, f, n, +, -\}, \delta, S_0, q_f)$

$\delta(S_0, c) = \{Q\};$

$\delta(Q, d) = \{Q\}; \delta(Q, k) = \{Q\}; \delta(Q, f) = \{Q\}; \delta(Q, +) = \{Q\}; \delta(Q, -) = \{Q\};$

$\delta(Q, n) = \{q_f\};$

Пример конфигурации КА: $S_0c+d-dkf-n \vdash Q+d-dkf-n \vdash Qd-dkf-n \vdash Q-dkf-n \vdash Qdkf-n \vdash Qkf-n \vdash Qf-n \vdash Q-n \vdash Qn \vdash q_f$

1.d. Определить свойства конечного автомата. Конечный автомат является детерминированным.

1.e. Построить диаграмму переходов КА рис. 2.6.:

{d, k, f, +, -}

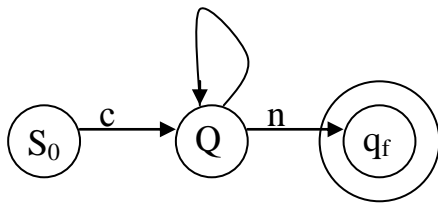


Рис. 2.6. Диаграмма переходов КА

2.6. Способы программной реализации

Способы программной реализации

Глава 3. Контекстно-свободные грамматики и МП-автоматы

3.1. Определение КС-грамматик и МП-автоматов

Из четырех типов грамматик иерархии Хомского класс контекстно-свободных грамматик наиболее важен с точки зрения приложений к языкам программирования и компиляции. С помощью этого типа грамматик определяется большая часть синтаксических структур языков программирования.

Определение 6. Контекстно-свободная грамматика типа 2:

Грамматика $G = (T, V, P, S_0)$ называется *контекстно-свободной* (КС) (или *бесконтекстной*), если каждое правило из P имеет вид $A \rightarrow \alpha$, где $A \in V$, $\alpha \in (T \cup V)^*$.

Заметим, что, согласно определению каждая праволинейная грамматика – КС. Языки, порождаемые КС-грамматиками, называются КС-языками.

Пример, КС-грамматики, порождающей скобочные арифметические выражения: $G_1 = (\{v, +, *, (,)\}, \{S, F, L\}, P, S)$, где $P = \{S \rightarrow S + F, S \rightarrow F, F \rightarrow F * L, F \rightarrow L, L \rightarrow v, L \rightarrow (S)\}$.

Пример вывода, заменяя самый левый нетерминал (левосторонний вывод) сентенциальной формы: $S \Rightarrow S + F \Rightarrow F + F \Rightarrow L + F \Rightarrow v + F \Rightarrow v + F * L \Rightarrow v + L * L \Rightarrow v + v * L \Rightarrow v + v * v$.

3.2. Преобразование КС-грамматик

Для построения синтаксических анализаторов необходимо, чтобы КС-грамматика была в *приведенной* форме см. рис. 3.1.

Если применить к КС-грамматике алгоритмы:

- устранения бесполезных символов

1). *Непроизводящих* и

2). *Недостижимых* символов;

3). преобразования в грамматику без ϵ -правил;

4). устранения цепных правил, то получим *приведенную* КС-грамматику.

Рассмотрим эти алгоритмы и их применение.

1. Символ $a \in T$ называется *недостижимым* в КС-грамматике G , если он не может появиться ни в одной из сентенциальных форм.

2. Нетерминальный символ $A \in V$ называется *производящим*, если из него можно вывести терминальную цепочку, т.е. если существует вывод $A \Rightarrow^+ \alpha$, где $\alpha \in T^+$. В противном случае символ называется *непроизводящим*.

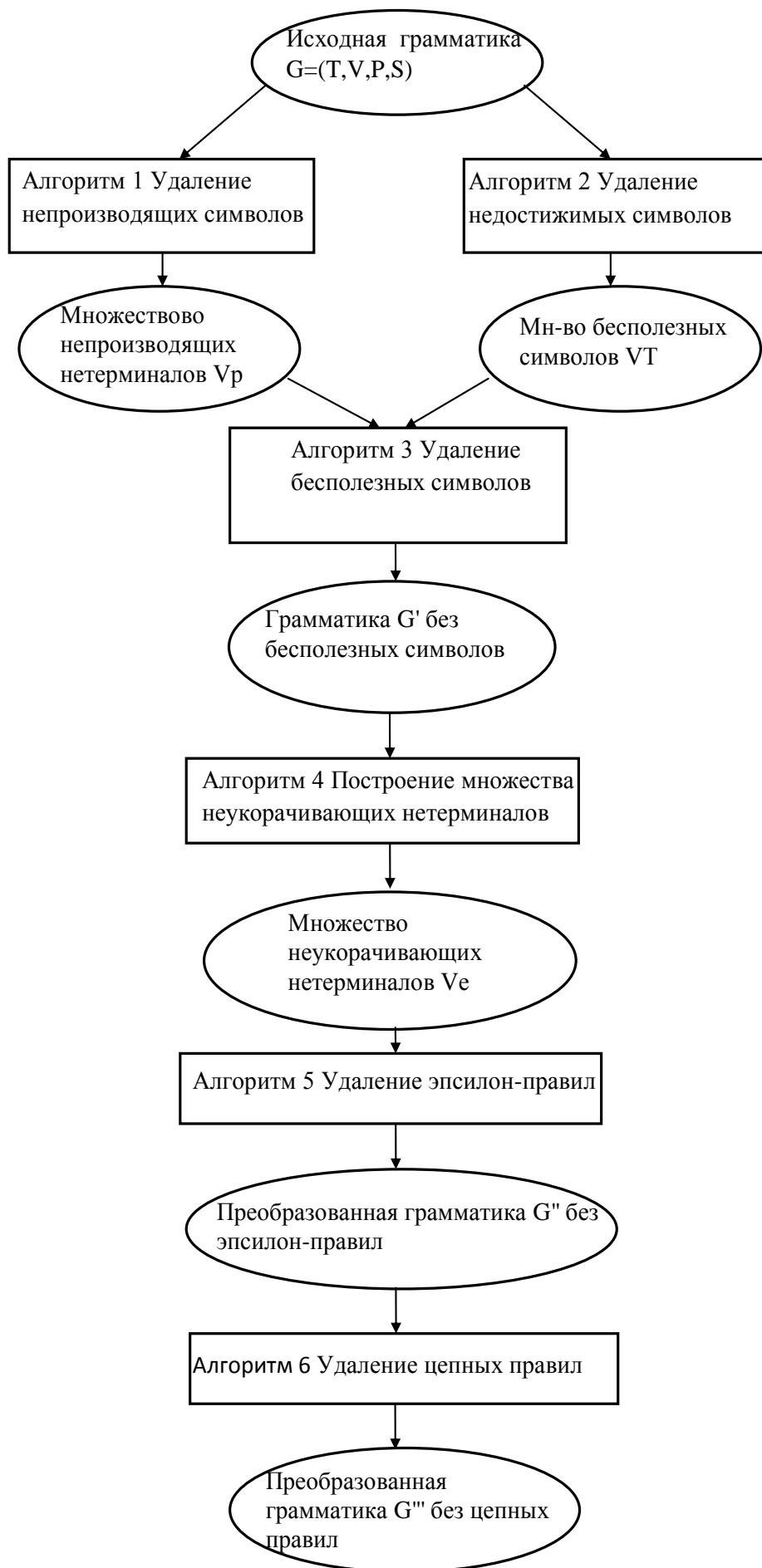


Рис.3.1 Метод преобразования грамматики к *приведенной* форме

1. Удаление непроеизводящих символов

Алгоритм 3.1. Определение множества *проеизводящих* нетерминальных символов V_p .

Если все символы цепочки из правой части правила вывода являются *проеизводящими*, то нетерминал в левой части правила вывода также должен быть *проеизводящим*.

Вход: КС $G = (T, V, P, S)$

Выход: $V_p = \{A \mid A \Rightarrow^+ \alpha, A \in V, \alpha \in T^+\}$

Алгоритм:

```
 $V_p := \emptyset;$  //множество проеизводящих символов пустое  
foreach (  $A \in V$  ) //проходим все нетерминалы  
  if  $((A \rightarrow \alpha X) \in P)$  //если есть проеизводящее правило  
     $V_p := V_p \cup A$  //добавляем нетерминал из левой части правила в множество  
  end if  
end foreach
```

Код на c#:

//определение множества *проеизводящих* нетерминальных символов

```
private ArrayList producingSymb() {  
    ArrayList Vp = new ArrayList();  
  
    foreach(Prule p in this.Prules) {  
        foreach(string t in this.T)  
            if (p.RightChain.Contains(t))  
                if (!Vp.Contains(p.leftNoTerm))  
                    Vp.Add(p.leftNoTerm); //если еще нет - добавляем  
    }  
    return Vp;  
}
```

2. Удаление недостижимых символов

Алгоритм 3.2. Определение множества *достижимых* символов V_r . Если нетерминал в левой части правила грамматики является *достижимым*, то *достижимы* и все символы правой части этого правила.

Вход: исходная КС $G = (T, V, P, S_0)$, при условии $\alpha, \beta \in (V \cup T)^*$, $X \in V$

Выход: VT_r

Алгоритм:

```
 $VT_r := \{S_0\}$   
foreach ( $A \in V$ ) //Проходим по всем нетерминалам  
  if  $(A \rightarrow \alpha X \beta) \in P$  //Если нетерминал достижим  
     $VT_r := VT_r \cup X$  //Добавляем его в множество  
    if  $\alpha \notin VT_r$  // Есть ли  $\alpha$  в  $VT_r$ ?  
       $VT_r := VT_r \cup \alpha$  // Если нет - добавляем  
    end if  
    if  $\beta \notin VT_r$ 
```

```

        VTr := VTr ∪ β
    end if
end if
end foreach

```

Код на с#:

```

//определение множества достижимых символов
private ArrayList Reachable(string StartState) {
    ArrayList Vr = new ArrayList(){ this.S0 };
//множество достижимых за 1 шаг
    ArrayList States = ReachableByOneStep(StartState);
//Нетерминальные символы из массива
    ArrayList NoTerm = NoTermReturn(States);
    Vr = Unify(Vr, NoTerm);
    foreach(string v in NoTerm) {
        Vr = Unify(Vr, ReachableByOneStep(v));
    }
    return Vr;
}

```

3. Удаление бесполезных символов

Алгоритм 3.3. Устранение бесполезных символов. Вначале исключить *непроизводящие* нетерминалы, а затем *недостижимые* символы.

Вход: исходная КС $G = (T, V, P, S)$, $V_p, VT_r, X \in V, \alpha \text{ и } \beta \in T^*$

Выход: КС $G' = (T', V', P', S)$

Алгоритм:

```

P' := ∅
foreach (A → αXβ) ∈ P // Проходим по всем правилам
    if X ∈ Vp //Если X – производящий нетерминал
        P' := P' ∪ (A → αXβ) //Добавляем правило, из которого выводится X
    end if
end foreach
T' := T ∩ VTr
V' := V ∩ VTr

```

Код на с#:

```

//удаление бесполезных символов
public myGrammar unUsefulDelete()
{
    ArrayList Vp = producingSymb(); //множество производящих символов
    ArrayList P1 = new ArrayList();
    ArrayList TorVp = Unify(Vp, this.T); //Объединение всех терминалов
    // и производящих символов
    foreach(Prule p in this.Prules)
    {
        ArrayList SymbInRule = SymbInRules(p); //все символы в правиле
        foreach(string s in SymbInRule)
            if (TorVp.Contains(s) || TermReturn(p.RightChain) !=
null) {

```

```

        if (!P1.Contains(p)) { P1.Add(p); }
    }
}
//построить множество достижимых символов
ArrayList Vr = Reachable(this.S0);
ArrayList T1 = intersection(this.T, Vr); //пересечение множеств
ArrayList V1 = intersection(Vr, this.V); //пересечение множеств
ArrayList P2 = new ArrayList();

foreach(Prule p in P1)
{
    // Все терминалы из правила
    ArrayList SymbInRule = SymbInRules(p);
    foreach(string s in SymbInRule) {
        if (Vr.Contains(s))
            if (!P2.Contains(p))
                P2.Add(p);
    }
}
return new myGrammar(T1, V1, P2, this.S0);
}

```

4. Построение множества неукорачивающих нетерминалов

Алгоритм 3.4. Построение множества V_ε укорачивающих нетерминалов.

Вход: КС $G' = (T, V, P, S)$.

Выход: V_ε

Алгоритм:

$V_\varepsilon := \emptyset$;

foreach A in V

if $(A \rightarrow \varepsilon) \in P$ //Если в правой части правила ε

$V_\varepsilon := V_\varepsilon \cup A$ //Добавляем его в множество правил

end if

end foreach

Код на c#:

//построение множества укорачивающих нетерминалов

```

private ArrayList ShortNoTerm() {
    ArrayList Ve = new ArrayList();
    foreach(Prule p in this.Prules) {
        if (p.rightChain.Contains("")) //если правые части содержат
пустые правила
            Ve.Add(p.leftNoTerm);
    }
    return Ve;
}

```

5. Удаление эpsilon-правил

Алгоритм 3.5. Преобразование КС-грамматики с ε -правилами в эквивалентную НКС-грамматику.

Вход: КС $G' = (T, V, P, S), V_\varepsilon$

Выход: НКС грамматика $G'' = (T, V', P', S')$

Алгоритм:

$P' := P$

foreach ($A \rightarrow \alpha_0 B_1 \alpha_1 B_2 \alpha_2 \dots B_k \alpha_k$) $\in P'$ // Проходим по всем правилам

// Если в правой части нетерминал, из которого выводится ε -правило

foreach ($B_k \in V_\varepsilon$)

$A_1 := \text{delete}(B_k \text{ from } A)$ // Кладем в A_1 правило A без ε -нетерминала

$P' := P' \cup A_1$ // Добавляем в P' правило A_1

end foreach

end foreach

if ($S \in V_\varepsilon$) //если нач. символ принадлежит множеству V_ε то добавляем новое правило

$V' = V \cup \{S'\}$

$P' := P' \cup (S' \rightarrow S \mid \varepsilon)$

end if

Код на с#:

//удаление эпсилон правил

```
public myGrammar EpsDelete() {
    Console.WriteLine("Deleting epsylon rules");

    ArrayList Ve = ShortNoTerm();
    ArrayList P1 = new ArrayList();
    ArrayList V1 = this.V;

    foreach(Prule p in Prules) {
        if (!ContainEps(p)) {
            P1.Add(p);
        }
    }

    // создаем новое правило
    Prule p1 = new Prule(p.leftNoTerm, TermReturn(p.rightChain));
    if (p.rightChain.Count != 1) {
        P1.Add(p1);
    }

    if (Ve.Contains(this.S0)) {
        V1.Add("S1");
        P1.Add(new Prule("S1", new ArrayList(){ this.S0 }));
        P1.Add(new Prule("S1", new ArrayList(){ "" }));
        return new myGrammar(this.T, V1, P1, "S1");
    }
    else
        return new myGrammar(this.T, V1, P1, this.S0);
}
```

6. Удаление цепных правил

Алгоритм 3.6. Устранение цепных правил.

Вход: КС $G'' = (T, V, P, S)$

Выход: НКС $G''' = (T, V, P', S)$

Алгоритм:

```
Pc := ∅
Pnc := ∅
P' := ∅

// нахождение цепных и нецепных правил
foreach (A → αBβ) ∈ P // Проходим по всем правилам
  if B ∈ V & α,β = ∅ // Если A → αBβ – цепное
    Pc := Pc ∪ (A → B) // Добавляем его в мн-во цепных Pc
  else
    Pnc := Pnc ∪ (A → αBβ) // Если нет, то добавляем в мн-во нецепных Pnc
  end if
end foreach

// замена цепных на нецепные
foreach (R → αBβ) ∈ P
  if B → γCδ ∈ Pc
    P' := P' ∪ (R → αγCδβ)
  endif
end foreach

P' := P' ∪ Pnc
```

Код на с#:

```
//удаление цепных правил
public myGrammar ChainRuleDelete() {
    Console.WriteLine("Цепные правила удалены...");

    ArrayList NoChainRules = new ArrayList();
    ArrayList ChainRules = new ArrayList();

    foreach (Prule p in this.Prules){
        if (TermReturn(p.rightChain) != null) { NoChainRules.Add(p); }
        else { ChainRules.Add(p); }
    }

    if (ChainRules.Count == 0) { Console.WriteLine("Цепных правил нет..."); }
    else
    {
        foreach (Prule chainrule in ChainRules)
            foreach (Prule rule in NoChainRules){
                if (rule.RightChain.Contains(chainrule.leftNoTerm))
                    for (int i = 0; i < rule.RightChain.Count; i++){
                        if (rule.RightChain[i].ToString() == chainrule.leftNoTerm)
                            rule.RightChain[i] = chainrule.RightChain[0].ToString();
                    }
            }
    }
    ArrayList P = new ArrayList();
    foreach (Prule rule in NoChainRules)
        if (!P.Contains(rule)) P.Add(rule);
    return new myGrammar(this.T, this.V, P, this.S0);
}
```

Пример применения метода приведения грамматики

Вход: КС грамматика $G (\{a, b, c, d\}, \{S, A, B, C, F\}, P, S)$, где P :

$S \rightarrow b \mid cAB$

$A \rightarrow Ab \mid c$

$B \rightarrow cB$

$C \rightarrow Ca$

$F \rightarrow d$

Результат работы программы:

Правила:

$S \rightarrow b$

$S \rightarrow A$

$S \rightarrow cAB$

$A \rightarrow Ab$

$A \rightarrow c$

$B \rightarrow cB$

$C \rightarrow Ca$

$F \rightarrow d$

Бесполезные символы удалены...

Правила:

$S \rightarrow b$

$S \rightarrow A$

$S \rightarrow cAB$

$A \rightarrow Ab$

$A \rightarrow c$

$B \rightarrow cB$

Эпсилон-правила удалены...

Правила:

$S \rightarrow b$

$S \rightarrow A$

$S \rightarrow cAB$

$S \rightarrow c$

$A \rightarrow Ab$

$A \rightarrow b$

$A \rightarrow c$

$B \rightarrow cB$

$B \rightarrow c$

Цепные правила удалены...

Правила:

$S \rightarrow b$

$S \rightarrow cAB$

$S \rightarrow c$

$A \rightarrow Ab$

$A \rightarrow b$

$A \rightarrow c$

$B \rightarrow cB$

$B \rightarrow c$

Левая рекурсия удалена...

$A \rightarrow Ab$

Правила:

$S \rightarrow b$

$S \rightarrow cAB$
 $S \rightarrow c$
 $A \rightarrow b$
 $A \rightarrow c$
 $B \rightarrow cB$
 $B \rightarrow c$
 $A \rightarrow S0$
 $S0 \rightarrow bS0$

Модификация

Данные алгоритмы можно модифицировать, объединяя похожие части. Самая простая модификация – объединение алгоритмов 1-3:

```

VTr := {S0}
Vp := ∅;
P' := ∅
foreach (A ∈ V)
  foreach (A → αXβ) ∈ P
    Vp := Vp ∪ A
    VTr := VTr ∪ X
    if α ∉ VTr
      VTr := VTr ∪ α
    end if
    if β ∉ VTr
      VTr := VTr ∪ β
    end if
    if X ∈ Vp
      P' := P' ∪ (A → αXβ)
    end if
  end foreach
end foreach
T' := T ∩ VTr
V' := V ∩ VTr

```

Таким же способом можно объединить и другие алгоритмы, но их сложность, как и в алгоритме выше, повысится, так как чтобы все алгоритмы работали в одном цикле, нужно будет создавать несколько дополнительных внутри него.

+++++

Строим рекурсивно множества $V_p^0, V_p^1, \dots, V_p^i, \dots$

1. Положить $V_p^0 = \emptyset, i=1$.
2. Положить $V_p^i = V_p^{i-1} \cup \{A \mid (A \rightarrow \alpha) \in P, \alpha \in T^+\}$.
3. Если $V_p^i \neq V_p^{i-1}$, положить $i = i + 1$ и перейти к шагу 2.
4. Положить $V_p = V_p^i$.

Поскольку $V_p \subseteq V$, то число повторений шага 2. $i+1$, где i – число нетерминальных символов грамматики G .

Алгоритм 3.2. Определение множества *достижимых* символов V_r . Если нетерминал в левой части правила грамматики является достижимым, то достижимы и все символы правой части этого правила.

Вход: КС грамматика $G = (T, V, P, S)$.

Выход: $V_r = \{X \mid S \Rightarrow^* \alpha X \beta, X \in V, \alpha, \beta \in (V \cup T)^*\}$.

Строим рекурсивно множества $V_r^0, V_r^1, \dots, V_r^i, \dots$

1. Положить $V_r^0 = \{S\}, i=1$.

2. Положить $V_r^i = V_r^{i-1} \cup \{X \mid A \rightarrow \alpha X \beta, X \in V, \alpha, \beta \in (V \cup T)^* \text{ и } A \in V_r^{i-1}\}$.

3. Если $V_r^i \neq V_r^{i-1}$, положить $i = i + 1$ и перейти к шагу 2.

4. Положить $V_r = V_r^i$.

3. Символ $a \in T \cup V$ называется *бесполезным* в КС-грамматике G , если он *непроизводящий* или *недостижимый*.

Алгоритм 3.3. Устранение *бесполезных* символов. Вначале исключить *непроизводящие* нетерминалы, а затем *недостижимые* символы.

Вход: КС грамматика $G = (T, V, P, S)$, для которой $L(G) \neq \emptyset$.

Выход: КС грамматика $G' = (T', V', P', S)$, у которой $L(G') \neq \emptyset$ и в $T' \cup V'$ нет бесполезных символов.

1. Построить множество V_p производящих нетерминалов грамматики G . Алгоритм 3.1.

2. Положить $G_1 = (T, V_p, P_1, S)$, где P_1 состоит из правил множества P , содержащих только символы из $T \cup V_p$ (алгоритмы 3.1.)

3. Построить множество V_r (алгоритм 3.2.), достижимых символов грамматики G_1 .

4. Положить $G' = (T', V', P', S)$, где $T' = T \cap V_r, V' = V_r \cap V, P'$ состоит из правил множества P_1 , содержащих только символы из множества V_r .

Алгоритмы 3.1 и 3.2 следует употреблять в указанном порядке.

Алгоритм устранения ε -правил в КС-грамматике основан на использовании множества укорачивающих нетерминалов.

КС-грамматика называется *неукорачивающей* КС-грамматикой (НКС-грамматикой, КС-грамматикой без ε -правил) при условии, что P не содержит $S \rightarrow \varepsilon$ и S не встречается в правах частях остальных правил.

Алгоритм 3.4. Построение множества V_ε укорачивающих нетерминалов.

Вход: КС грамматика $G = (T, V, P, S)$.

Выход: $V_\varepsilon = \{A \mid A \Rightarrow^+ \varepsilon, A \in V\}$.

Строим рекурсивно множества укорачивающих нетерминалов $V_\varepsilon^0, V_\varepsilon^1, \dots, V_\varepsilon^i, \dots$

1. Положить $V_\varepsilon^0 = \{A \mid A \rightarrow \varepsilon, A \in V\}, i=1$.

2. Положить $V_\varepsilon^i = V_\varepsilon^{i-1} \cup \{A \mid (A \rightarrow \alpha) \in P \text{ и } \alpha \in (V_\varepsilon^{i-1})^+\}$.

3. Если $V_\varepsilon^i \neq V_\varepsilon^{i-1}$, положить $i = i + 1$ и перейти к шагу 2.

4. Положить $V_\varepsilon = V_\varepsilon^i$.

Поскольку $V_\varepsilon \subseteq V$, то число повторений шага 2. $i+1$, где i – число нетерминальных символов грамматики G .

Алгоритм 3.5. Преобразование КС-грамматики с ε -правилами в эквивалентную НКС-грамматику.

Вход: КС грамматика $G = (T, V, P, S)$.

Выход: НКС грамматика $G' = (T, V', P', S')$, порождающая язык $L(G') = L(G)$.

1. Построить множество V_ε , укорачивающих нетерминалов (алгоритмы 3.4.).
2. Положить $P' = \emptyset$.
3. Если $(A \rightarrow \alpha_0 B_1 \alpha_1 B_2 \alpha_2 \dots B_k \alpha_k) \in P$, где $k \geq 0$, $B_j \in V_\varepsilon$ ($1 \leq j \leq k$) и ни один из символов цепочек $\alpha_i \in (V \cup T)^*$ ($0 \leq i \leq k$) не содержит символов из V_ε , то включить в P' , все правила вывода вида $A \rightarrow \alpha_0 X_1 \alpha_1 X_2 \alpha_2 \dots X_k \alpha_k$, где $X_j = B_j$ или $X = \varepsilon$.
Если все цепочки $\alpha_i = \varepsilon$, то правила $A \rightarrow \varepsilon$ не включать в P' .
4. Если $S \in V_\varepsilon$, то положить $V' = V \cup \{S'\}$ и добавить в P' два правила: $S' \rightarrow S$ и $S' \rightarrow \varepsilon$. В противном случае положить $V' = V$ и $S' = S$.
5. Положить $G' = (T, V', P', S')$.

В КС грамматиках часто встречаются правила, правая часть которых состоит из одного нетерминального символа. Правило вида $A \rightarrow B$ называется цепным правилом.

Алгоритм 3.6. Устранение цепных правил.

Вход: КС грамматика $G = (T, V, P, S)$.

Выход: Эквивалентная НКС грамматика $G' = (T, V, P', S)$ без цепных правил.

1. Для каждого $A \in V$, построить множество $V_A = \{B \mid A \Rightarrow^+ B\}$ следующим образом:
 - 1.1. Положить $V_A^0 = \{A\}$, $i=1$.
 - 1.2. Положить $V_A^i = V_A^{i-1} \cup \{C \mid (B \rightarrow C) \in P \text{ и } B \in V_A^{i-1}\}$.
 - 1.3. Если $V_A^i \neq V_A^{i-1}$, положить $V_A = V_A^i$.
2. Положить $P' = \emptyset$.
3. Если $B \rightarrow \alpha \in P$ и не является цепным правилом, то положить $P' = P' \cup \{A \rightarrow \alpha\}$ для всех таких A , что $B \in V_A$.
4. Положить $G' = (T, V, P', S)$.

Нетерминал КС-грамматики называется *рекурсивным*, если $A \Rightarrow^+ \alpha A \beta$, для некоторых α и β . Если $\alpha = \varepsilon$, то A называется *леворекурсивным*, если $\beta = \varepsilon$, то A называется *праворекурсивным*. Грамматика, имеющая хотя бы один леворекурсивный нетерминал, называется *леворекурсивной*. Грамматика, имеющая хотя бы один праворекурсивный, нетерминал называется *праворекурсивной*.

При нисходящем синтаксическом анализе требуется, чтобы приведенная грамматика рассматриваемого языка не содержала **левой рекурсии**.

Для любой КС-грамматики существует эквивалентная грамматика без левой рекурсии. Рассмотрим алгоритм устранения левой рекурсии.

Алгоритм 3.7.

Вход: Приведенная КС грамматика $G = (T, V, P, S_0)$.

Выход: Эквивалентная КС грамматика G' .

1. Пусть $V = \{A_1, \dots, A_n\}$. Преобразуем G так, чтобы в правиле $A_i \rightarrow \alpha$, цепочка α начиналась либо с терминала, либо с такого A_j , что $j > i$. Пусть $i=1$.
 2. Пусть множество A_i правил – это $A_i \rightarrow A_i\alpha_1 \mid \dots \mid A_i\alpha_m \mid \beta_1 \mid \dots \mid \beta_p$, где ни одна цепочка β_j не начинается с A_k , если $k \leq i$. Заменяем A_i -правила правилами:

$$A_i \rightarrow \beta_1 \mid \dots \mid \beta_p \mid \beta_1 A_i' \mid \dots \mid \beta_p A_i'$$

$$A_i \rightarrow \alpha_1 \mid \dots \mid \alpha_m \mid \alpha_1 A_i' \mid \dots \mid \alpha_m A_i'$$
где A_i' – новый символ. Правые части всех A_i -правил начинаются теперь с терминала или с A_k для некоторого $k > i$.
 3. Если $i = n$, то останов и получена грамматика G' , иначе $j = i, i = i + 1$.
 4. Заменить каждое правило вида $A_i \rightarrow A_j\alpha$ правилами $A_i \rightarrow \beta_1\alpha \mid \dots \mid \beta_m\alpha$, где $A_j \rightarrow \beta_1 \mid \dots \mid \beta_m$ – все A_j -правила.
- Так как правая часть каждого A_j – правила начинается уже с терминала или с A_k для $k > j$, то и правая часть каждого A_i – правила будет обладать этим же свойством.
5. Если $j = i - 1$, перейти к шагу 2, иначе $j = j + 1$ и перейти к шагу 4.

Частный случай не леворекурсивной грамматики – грамматика в нормальной форме Шейлы Грейбах.

Определение 7. КС грамматика $G = (T, V, P, S)$ называется грамматикой в нормальной форме Грейбах, если в ней нет ε -правил, т.е. правил вида $A \rightarrow \varepsilon$, и каждое правило из P отличное от $S \rightarrow \varepsilon$, имеет вид $A \rightarrow a\alpha$, где $a \in T, \alpha \in V^*$.

Также полезно представлять грамматику в нормальной форме Хомского, что позволяет упростить рассмотрение ее свойств.

Определение 8. КС грамматика $G = (T, V, P, S)$ называется грамматикой в нормальной форме Хомского, если каждое правило из P имеет один из следующих видов:

1. $A \rightarrow BC$, где $A, B, C \in V$;
2. $A \rightarrow a$, где $a \in T$;
3. $S \rightarrow \varepsilon$, если $\varepsilon \in L(G)$, причем S не встречается в правых частях правил.

3.3. Пример: приведенная форма КС-грамматики

Постановка задачи. Привести КС-грамматику к *приведенной* форме, исключить левую рекурсию.

А). Устранить из грамматики G бесполезные символы. Применить алгоритм 3.3. к грамматике $G = (T, V, P, S)$, где $V = \{A, B, C, S\}$, $T = \{a, b, c\}$, $P = \{S \rightarrow aC, S \rightarrow A, A \rightarrow cAB, B \rightarrow b, C \rightarrow a\}$.

Шаг 1. множество $V_p = \{B, C, S\}$ A – *непроизводящий символ*

Шаг 2. $G_1 = (\{B, C, S\}, \{a, b, c\}, P, S)$, где $P = \{S \rightarrow aC, B \rightarrow b, C \rightarrow a\}$,

Шаг 3. $V_r = \{C, S\}$ B – *недостижим, так как A-непроизводящий символ*

Шаг 4. $G' = (\{C, S\}, \{a\}, P, S)$, где $P = \{S \rightarrow aC, C \rightarrow a\}$, $L(G') = \{aa\}$

В). Устранить из грамматики G ε -правила, применить алгоритм 3.5. Преобразовать грамматику $G = (T, V, P, S)$, где $V = \{A, S\}$, $T = \{b, c\}$, $P = \{S \rightarrow cA, S \rightarrow \varepsilon, A \rightarrow cA, A \rightarrow bA, A \rightarrow \varepsilon\}$, в эквивалентную НКС-грамматику.

Шаг 1. Применяя алгоритм 3.4., получаем $V_\varepsilon^0 = \{S, A\}$, $V_\varepsilon^1 = \{S, A\}$, значит $V_\varepsilon^0 = V_\varepsilon^1 = \{S, A\}$.

Шаг 2. Положить $P' = \emptyset$.

Шаг 3. Алгоритм 3.5. Рассмотрим правило $S \rightarrow cA$.

Для него в новое множество правил грамматики P' добавляем правила $S \rightarrow cA$ и $S \rightarrow c$.

Для правила $A \rightarrow cA$, добавляем в P' правила $A \rightarrow cA$ и $A \rightarrow c$.

Для правила $A \rightarrow bA$, добавляем в P' $A \rightarrow bA$ и $A \rightarrow b$,

тогда $P' = \{S \rightarrow cA, S \rightarrow c, A \rightarrow cA, A \rightarrow c, A \rightarrow bA, A \rightarrow b\}$.

Шаг 4. $S \in V_\epsilon$, то в V' добавляем новый нетерминал S' , а в P' два правила $S' \rightarrow S$ и $S' \rightarrow \epsilon$.

Шаг 5. $G' = (\{b, c\}, \{S', S, A\}, \{S' \rightarrow S, S' \rightarrow \epsilon, S \rightarrow cA, S \rightarrow c, A \rightarrow cA, A \rightarrow c, A \rightarrow bA, A \rightarrow b\}, S')$.

С). Устранить из КС грамматики G цепные правила, применить алгоритм 3.6. $G = (\{S, F, L\}, \{v, +, *, (,)\}, P, S)$, где P состоит из правил

$S \rightarrow S + F \mid F, F \rightarrow F * L \mid L, L \rightarrow v \mid (S)$

Шаг 1. Применяя алгоритм 3.6., получаем $V_S = \{S, F, L\}$, $V_F = \{F, L\}$, $V_L = \{L\}$.

Шаг 2. Положить $P' = \emptyset$.

Шаг 3. Выберем первый нетерминал S из множества V_S . Множество правил, левая часть которых – нетерминальный символ S , правые части – это правые части нецепных правил исходной грамматики, в левой части которых находятся символы из множества V_S , получаем: $\{S \rightarrow S + F \mid F * L \mid (S) \mid v\}$. Включаем эти правила в $P' = \{S \rightarrow S + F \mid F * L \mid (S) \mid v\}$. Таким же образом рассматриваем символы из V_F, V_L ,

Шаг 4. В результате $G' = (T, V, P', S)$, где $P' = \{S \rightarrow S + F \mid F * L \mid (S) \mid v, F \rightarrow F * L \mid (S) \mid v, L \rightarrow v \mid (S)\}$

Д). Исключить левую рекурсию из КС – грамматики G .

1. Устранить левую рекурсию в заданной КС-грамматике G , порождающей скобочные арифметические выражения.

$G = (\{v, +, *, (,)\}, \{S, F, L\}, P, S)$, где P состоит из правил

$S \rightarrow S + F \mid F$

$F \rightarrow F * L \mid L$

$L \rightarrow v \mid (S)$

Применить алгоритм 3.7. к грамматике G . Указаны шаги алгоритма 3.7.

Шаг 1. Пусть $A_1 = S, A_2 = F, A_3 = L$. $V = \{A_1, A_2, A_3\}$.

Шаг 2. Для $i = 1$ преобразуем правила: $S \rightarrow S + F \mid F$, $\alpha = + F$, $\beta = F$. Заменим S - правила правилами: $S \rightarrow F \mid FS'$, и добавим в грамматику правило для нового нетерминала $S' \rightarrow + F \mid + FS'$.

Шаг 3. $i = 2$ и $j = 1$.

Шаг 4. Для $i = 2, j = 1$ правила вида $A_i \rightarrow A_j \alpha$ отсутствуют.

Шаг 5. Так как $j = i - 1$, то переходим к **шагу 2**.

Шаг 2. Для $i = 2$ преобразуем правила: $F \rightarrow F * L \mid L$, $\alpha = * L$, $\beta = L$. Заменим F - правила правилами: $F \rightarrow L \mid LF'$, и добавим в грамматику правило для нового нетерминала $F' \rightarrow * L \mid LF'$.

Шаг 3. $i = 3, j = 1$.

Шаг 4. Для $i = 3, j = 1$ правила вида $A_i \rightarrow A_j \alpha$ отсутствуют.

Шаг 5. Так как $j \neq i - 1$, то $j = j + 1$ переходим к **шагу 4**.

Шаг 4. Для $i = 3, j = 2$ правила вида $A_i \rightarrow A_j \alpha$ отсутствуют.

Шаг 5. Так как $j = i - 1$, то переходим к **шагу 2**.

Шаг 2. Для $i = 3$ правила вида $A_i \rightarrow A_i \alpha$ отсутствуют.

Шаг 3. $i = n = 3$, преобразование завершено.

Получили правила новой грамматики G' :

$S \rightarrow F \mid FS', S' \rightarrow +F \mid +FS', F \rightarrow L \mid LF', F' \rightarrow *L \mid LF', L \rightarrow (S + F) \mid v \mid (S).$

F) Определить в какой форме (Грейбах, Хомского) находится КС-грамматика G' .

Е). G' – приведенная КС-грамматика.

3.4. Определение МП-автоматов

Автоматы с магазинной памятью (МП - автоматы) представляют собой модель распознавателей для языков, задаваемых КС-грамматиками. МП - автоматы имеют вспомогательную память, называемую магазином см. 1.10. В магазин можно поместить неограниченное количество символов. В каждый момент времени доступен только верхний символ магазина.

Верхний символом магазина будем считать самый левый символ цепочки.



Рис. 1.10. Модель МП - автомата

Определение 9. МП автомат – это семерка объектов

$МП = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$

Q – конечное множество состояний устройства управления;

Σ – конечный алфавит входных символов;

Γ – конечный алфавит магазинных символов;

δ – функция переходов, отображает множества $Q \times (\Sigma \cup \{\varepsilon\}) \times \Gamma$ в множество конечных подмножеств множества $Q \times \Gamma^*$;

q_0 – начальное состояние, $q_0 \in Q$;

z_0 – начальный символ магазина, $z_0 \in \Gamma$;

F – множество заключительных состояний, $F \subseteq Q$.

Определение 10. Конфигурацией МП-автомата называется тройка $(q, \omega, z) \in Q \times \Sigma^* \times \Gamma^*$, где

q – текущее состояние управляющего устройства;

ω – необработанная часть входной цепочки (первый символ цепочки ω находится под входной головкой; если $\omega = \varepsilon$, то считается, что вся входная цепочка прочитана);

z – содержимое магазина (самый левый символ цепочки z считается верхним символом магазина; если $z = \varepsilon$, то магазин считается пустым).

Такт МП-автомата будем описывать бинарным отношением \vdash , определенным на множестве конфигураций. Будем писать:

$$(q, a\omega, z) \vdash (q', \omega, z\gamma), \text{ если } \delta(q, a, z) = (q', \gamma), \text{ где} \\ q, q' \in Q, a \in \Sigma \cup \{\varepsilon\}, \omega \in \Sigma^*, z \in \Gamma \text{ и } \gamma \in \Gamma^*$$

Если $a \neq \varepsilon$, то и входная цепочка прочитана не вся, то запись $(q, a\omega, z\gamma) \vdash (q', \omega, a\gamma)$ означает, что МП-автомат в состоянии q , обозревая символ во входной цепочки и имея символ z в верхушке магазина, может перейти в новое состояние q' , сдвинуть входную головку на один символ вправо и заменить верхний символ магазина z цепочкой магазинных символов γ .

Если $z = a$, то верхний символ удаляется из магазина.

Если $a = \varepsilon$, то текущий входной символ в этом такте называется ε -тактом, не принимается во внимание и входная головка остается неподвижной.

ε -такты могут выполняться также в случае, когда вся входная цепочка прочитана, но если магазин пуст, то такт МП-автомата невозможен по определению.

Так же, как и для конечных автоматов, можно определить транзитивное \vdash^+ и рефлексивно-транзитивное \vdash^* замыкания (отношения \vdash).

Начальной конфигурацией МП-автомата называется конфигурация вида (q_0, ω, z_0) , где устройство управления находится в начальном состоянии, на входной ленте записана цепочка $\omega \in \Sigma^*$, которую необходимо распознать, а магазин содержит только начальный символ z_0 .

Заключительной конфигурацией МП-автомата называется конфигурация вида (q, ε, γ) , где $q \in F$ – одно из заключительных состояний устройства управления, входная цепочка прочитана до конца, а в магазине записана некоторая, заранее определенная цепочка $\gamma \in \Gamma^*$.

Есть два способа определить язык, допускаемый МП-автоматом:

1. множеством входных цепочек, для которых существует опустошающая магазин последовательность операций;
2. множеством входных цепочек, для которых существует последовательность операций, приводящая автомат в заключительное состояние.

Цепочка $\omega \in \Sigma^*$ допускается МП-автоматом $МП = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$, если $(q_0, \omega, z_0) \vdash^* (q, \varepsilon, \gamma)$ для некоторых $q \in F$ и $\gamma \in \Gamma^*$.

Язык распознаваемый МП-автоматом, называется множество цепочек:

$$L(МП) = \{\omega \mid \omega \in \Sigma^* \text{ и } (q_0, \omega, z_0) \vdash^* (q, \varepsilon, \gamma) \text{ для некоторых } q \in F \text{ и } \gamma \in \Gamma^*\}$$

Пример. Определим МП-автомат, допускающий язык $L = \{a^n b^n \mid n \geq 0\}$
Цепочка символов языка L : aaabbb

$$МП = (\{q_0, q_1, q_2, q_f\}, \{a, b\}, \{z_0, a\}, \delta, q_0, z_0, \{q_f\})$$

$$\delta(q_0, a, z_0) = \{(q_1, az_0)\}$$

$$\delta(q_1, a, a) = \{(q_1, aa)\}$$

$$\delta(q_1, b, a) = \{(q_2, \varepsilon)\}$$

$$\delta(q_2, b, a) = \{(q_2, \varepsilon)\}$$

$$\delta(q_2, \varepsilon, z_0) = \{(q_f, \varepsilon)\}$$

Последовательность тактов:

$$\begin{aligned} (q_0, aaabbb, z_0) &\vdash_1 (q_1, aabbb, az_0) \vdash_2 (q_1, abbb, aaz_0) \vdash_3 \\ (q_1, bbb, aaaz_0) &\vdash_3 (q_2, bb, aaz_0) \vdash_4 (q_2, b, az_0) \vdash_4 \\ (q_2, \varepsilon, z_0) &\vdash_5 (q_f, \varepsilon, \varepsilon) \end{aligned}$$

Алгоритм 3.8. По КС-грамматике $G = (T, V, P, S)$ можно построить МП-автомат, $L(\text{МП}) = L(G)$. Пусть $\text{МП} = (\{q\}, \Sigma, \Sigma \cup V, \delta, q, S, \{q\})$, где δ определяется следующим образом:

1. Если $A \rightarrow \alpha$ - правило грамматики G , то $\delta(q, \varepsilon, A) = (q, \alpha)$.
2. $\delta(q, a, a) = \{(q, \varepsilon)\}$ для всех $a \in \Sigma$.

Пример проектирования МП-автомата. Постановка задачи. Построить МП-автомат и расширенный МП-автомат по КС-грамматике $G = (T, V, P, S)$, без левой рекурсии, автомат распознает скобочные выражения. Написать последовательность тактов для выделенной цепочки. Определить свойства автомата.

1. А). Построить МП-автомат, распознающий скобочные арифметические выражения заданные КС-грамматикой (не приведенная) $G_1 = (\{v, +, *, (,)\}, \{S, F, L\}, P, S)$, где P состоит из правил: $S \rightarrow S + F \mid F, F \rightarrow F * L \mid L, L \rightarrow v \mid (S)$.

Используя алгоритм 3.8. получим: $\text{МП} = (\{q\}, \{v, +, *, (,)\}, \{v, +, *, (,), S, F, L\}, \delta, q_0, S, \{q\})$, в котором функция переходов δ определяется следующим образом:

1. $\delta(q_0, \varepsilon, S) = \{(q, S + F), (q, F)\}$;
2. $\delta(q, \varepsilon, F) = \{(q, F * L), (q, L)\}$;
3. $\delta(q, \varepsilon, L) = \{(q, v), (q, (S))\}$;
4. $\delta(q, a, a) = \{(q, \varepsilon)\}$ для всех $a \in \Sigma = \{v, +, *, (,)\}$.

В). Последовательность тактов МП-автомата для цепочки $v^*(v+v)$:

$$\begin{aligned} (q_0, v * (v + v), S) &\vdash_1 (q, v * (v + v), F) \vdash_2 (q, v * (v + v), F * L) \vdash_3 (q, v * (v + v), \\ L * L) &\vdash_4 (q, v * (v + v), v * L) \vdash_5 (q, *, (v + v), * L) \vdash_6 (q, (v + v), L) \vdash_7 (q, (v + v), (S)) \\ &\vdash_8 (q, v + v, S)) \vdash_9 (q, v + v, S + F)) \vdash_{10} (q, v + v, F + F)) \vdash_{11} (q, v + v, L + F)) \vdash_{12} (q, \\ v + v, v + F)) &\vdash_{13} (q, +v, +F)) \vdash_{14} (q, v, F)) \vdash_{15} (q, v, L)) \vdash_{16} (q, v, v)) \vdash_{17} (q,),)) \vdash_{18} \\ (q, \varepsilon, \varepsilon). \end{aligned}$$

Последовательность тактов, которую выполняет МП-автомат, соответствует левому выводу цепочки $v * (v+v)$ в грамматике G_1 .

Тип синтаксических анализаторов, которые можно построить таким образом называют *нисходящим* (предсказывающим) анализатором.

Синтаксические анализаторы, построенные на основе расширенного МП-автомата, называют *восходящими* анализаторами, вывод строится снизу в верх.

Определение 11. Расширенным РМП-автоматом называется семерка объектов $\text{РМП} = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$, где верхним элементом магазина является самый правый символ цепочки, δ -функция переходов, которая задает

отображение множества $Q \times (\Sigma \cup \{\varepsilon\}) \times \Gamma$ в множество конечных подмножеств множества $Q \times \Gamma^*$, а все остальные объекты такие же, как и у МП-автомата см. определение 9.

Конфигурация расширенного МП-автомата определяется так же, как и для МП-автомата.

Алгоритм 3.9. По КС-грамматике $G = (T, V, P, S)$ можно построить расширенный РМП автомат, $L(\text{РМП}) = L(G)$. Будем обозначать символом \perp - конец цепочки. Пусть $\text{РМП} = (\{q, r\}, \Sigma, \Sigma \cup V \cup \perp, \delta, q, \perp, \{r\})$, где δ определяется следующим образом:

1. $\delta(q, a, \varepsilon) = \{(q, a)\}$ для всех $a \in \Sigma$ (символы с входной ленты на этих тактах переносятся в магазин).
2. Если $A \rightarrow \alpha$ - правило вывода грамматики G , то $\delta(q, \varepsilon, \alpha) = (q, A)$.
3. $\delta(q, \varepsilon, \perp S) = \{(r, \varepsilon)\}$

РМП-автомат продолжает работу пока магазин не станет пустым.

На практике часто используются детерминированные МП-автоматы.

Пример построения расширенного РМП-автомата

А). Построение расширенного РМП-автомата $(Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$, используем алгоритм 3.9, получим $\text{РМП} = (\{q, r\}, \{v, +, *, (,)\}, \{v, +, *, (,), S, F, L, \perp\}, \delta, q, \perp, \{r\})$, где функция переходов δ определена следующим образом:

1. $\delta(q, a, \varepsilon) = \{(q, a)\}$ для всех $a \in \Sigma = \{v, +, *, (,)\}$;
2. $\delta(q, \varepsilon, S+F) = \{(q, S)\}$;
3. $\delta(q, \varepsilon, F) = \{(q, S)\}$;
4. $\delta(q, \varepsilon, F * L) = \{(q, F)\}$;
5. $\delta(q, \varepsilon, L) = \{(q, F)\}$;
6. $\delta(q, \varepsilon, v) = \{(q, L)\}$;
7. $\delta(q, \varepsilon, (S)) = \{(q, L)\}$;
8. $\delta(q, \varepsilon, \perp S) = \{(r, \varepsilon)\}$;

В). При анализе входной цепочки $v+v*v$ расширенный РМП-автомат выполнит следующую последовательность тактов:

$(q, v+v*v, \perp) \xrightarrow{1} (q, +v*v, \perp v)$
 $\xrightarrow{6} (q, +v*v, \perp L)$
 $\xrightarrow{5} (q, +v*v, \perp F)$
 $\xrightarrow{3} (q, +v*v, \perp S)$
 $\xrightarrow{1} (q, v*v, \perp S+)$
 $\xrightarrow{1} (q, *v, \perp S+v)$
 $\xrightarrow{6} (q, *v, \perp S+L)$
 $\xrightarrow{5} (q, *v, \perp S+F)$
 $\xrightarrow{1} (q, v, \perp S+F*)$
 $\xrightarrow{1} (q, \varepsilon, \perp S+F*v)$
 $\xrightarrow{6} (q, \varepsilon, \perp S+F*L)$
 $\xrightarrow{4} (q, \varepsilon, \perp S+F)$
 $\xrightarrow{2} (q, \varepsilon, \perp S)$
 $\xrightarrow{8} (r, \varepsilon, \varepsilon)$

МП-автомат и расширенный РМП-автоматы – детерминированные автоматы.

Обычно синтаксический анализ выполняется путем моделирования МП-автомата, анализирующего входные цепочки.

МП-автомат отображает входные цепочки в соответствующие левые выводы (нисходящий анализ) или правые разборы (восходящий анализ).

Пусть задана КС-грамматика $G = (T, V, P, S)$, правила которой пронумерованы числами $1, 2, \dots, p$, и цепочка $\alpha \in (T \cup V)^*$,

- левым разбором цепочки α называется последовательность правил, примененных при ее левом выводе из S ;
- правым разбором цепочки α называется обращение последовательности правил, примененных при ее левом выводе из S .

Определение 12. Детерминированным МП-автоматом (ДМП-автоматом) с магазинной памятью называется МП-автомат $ДМП = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$, у которого для каждого $q \in Q$ и $z \in \Gamma$ выполняется одно из условий:

1. $\delta(q, a, z)$ содержит не более одного элемента для каждого $a \in \Sigma$ и $\delta(q, \varepsilon, z) = \emptyset$;
2. $\delta(q, a, z) = \emptyset$ для всех $a \in \Sigma$, и $\delta(q, \varepsilon, z)$ содержит не более одного элемента.

Если каждое правило КС-грамматики начинается с терминального символа, причем альтернативы начинаются с различных символов, то достаточно просто построить *детерминированный* синтаксический анализатор.

3.5. Способы реализации синтаксических анализаторов

Рассмотрим три способа реализации синтаксических анализаторов для заданного языка $L(G)$. Каждый имеет свои достоинства и недостатки. Напомним, что G является приведенной.

1. МП-автомат строим для $L(G) = L(MP)$, правила автомата строятся в соответствии с рассмотренными алгоритмами.

2. Строятся диаграммы Вирта, для каждой грамматики, создается неявный стек при косвешней рекурсии, что не требует реализации.

3. Таблично-управляемый разбор (предварительно строится диаграмма Вирта), диаграмма реализуется в виде списочной структуры, над которой выполняется алгоритм.

Диаграммы Вирта

1. Постановка задачи. Реализовать синтаксический анализатор левым разбором, используя приведенную грамматику без левой рекурсии $G = (\{a, b, c, -, +, (\), \}, \{E, T, F, \text{or}\}, P, S)$, где $P = \{E \rightarrow T \text{ or } T, \text{ or} \rightarrow + \mid -, T \rightarrow F \mid (E), F \rightarrow a \mid b \mid c\}$

2. Шаги реализации.

Выделить цепочку, принадлежащую языку задаваемому КС-грамматикой G , например: $L(G) = a + (b - c)$.

Привести вывод выделенной цепочки:

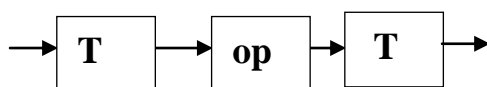
$E \Rightarrow T \text{ or } T \Rightarrow T \text{ or } (E) \Rightarrow F \text{ or } (T \text{ or } T) \Rightarrow a \text{ or } (T \text{ or } T) \Rightarrow a + (T \text{ or } T) \Rightarrow a + (F \text{ or } T) \Rightarrow a + (b \text{ or } T) \Rightarrow a + (b \text{ or } F) \Rightarrow a + (b \text{ or } c) \Rightarrow a + (b - c)$

Пронумеровать правила грамматики и представить их, используя метасимволы $\{\dots\}$, $:=$ расширенной формы Бекуса-Наура:

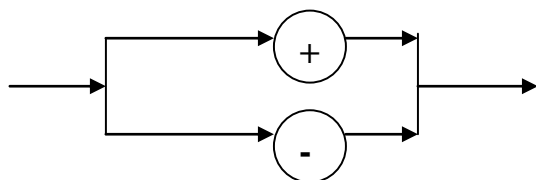
$p_1: E := \{ T \text{ or } T \}$, $p_2: \text{or} := \{ +, - \}$, $p_3: T := \{ F, (E) \}$, $p_4: F := \{ a, b, c \}$

Составить синтаксический граф, диаграмму Вирта для детерминированного грамматического разбора с просмотром вперед на один символ. Никакие две ветви не должны начинаться с одного и того же символа, если какой-либо граф А можно пройти, не читая вообще никаких входных символов, то такая “нулевая ветвь” должна помечаться всеми символами, которые могут следовать за А (это необходимо для принятия решения о переходе на эту ветвь).

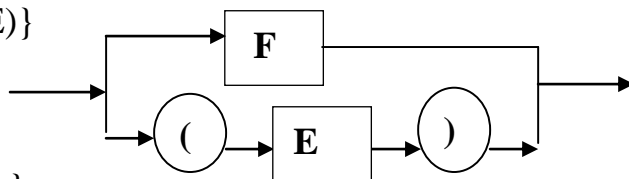
$p_1: E := \{ T \text{ or } T \}$



$p_2: \text{or} := \{ +, - \}$



$p_3: T := \{ F, (E) \}$



$p_4: F := \{ a, b, c \}$

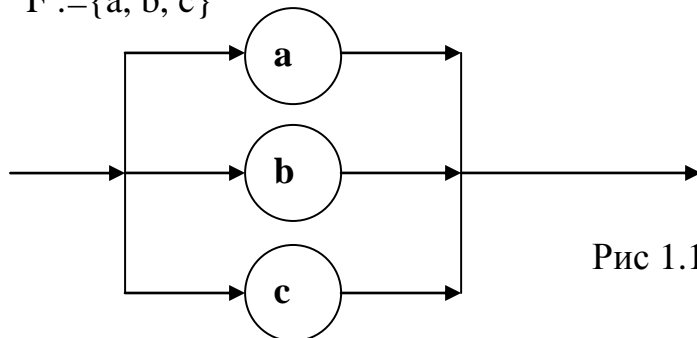
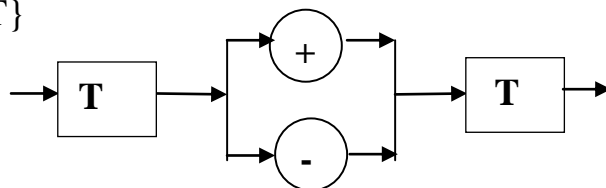


Рис 1.11. Диаграммы Вирта.

Правила преобразования графа в программу:

1. Свести систему графов к возможно меньшему числу отдельных графов с помощью соответствующих подстановок. Сделаем две подстановки, p_2 подставим в p_1 , p_4 подставим в p_3 :

$p_1: E := \{ T \{ +, - \} T \}$



$p_4: T := \{ \{a, b, c\}, (E) \}$

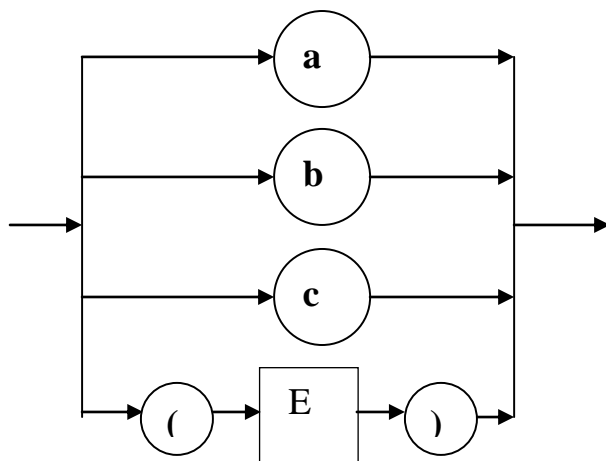


Рис 1.12. Преобразованные диаграммы Вирта.

Полученные правила на графе соответствуют нормальной форме Грейбах.

Преобразовать каждый синтаксический граф в описание процедуры в соответствии с правилами:

1. Каждому графу (диаграмме Вирта) ставится в соответствие процедура.
2. Элемент графа, обозначающий другой граф, переводится в оператор обращения к процедуре.
3. Последовательность элементов в последовательный вызов операторов. Выбор элементов в switch или условный оператор if.

Входные и выходные данные. Пользователь вводит строку в поле редактирования и получает ответ: Да – строка принадлежит языку, порожаемому данной грамматикой, и Нет – строка не принадлежит этому языку. Программа распознает язык $L(G)$ по заданной КС-грамматике G .

Таблично-управляемый разбор

1. Грамматика задается в виде данных, процедура разбора универсально для всех грамматик. Программа работает в строгом соответствии с методом простого нисходящего грамматического разбора и основывается на детерминированном синтаксическом графе (т.е. предложение должно анализироваться просмотром вперед на один символ без возврата).

2. Естественный способ представить граф в виде структуры данных, а не программ – это ввести узел для каждого символа и связать эти узлы с помощью ссылок. Выделим два типа узлов для терминальных (идентифицируется терминальным символом) и нетерминальных символов (ссылка на данные нетерминального символа). Тогда структура узла графически может быть представлена как:

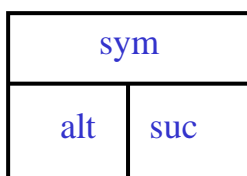


Рис. 1. 14. Структура узла Node

Где “suc” - последователь, “alt” – альтернатива. Пустую последовательность обозначим “empty”. Конструируется класс Node (Узел) для данной структуры.

3. Правила преобразования графов, в структуру данных: получить как можно меньшее число графов с помощью подстановок.

4. Последовательность элементов преобразовать в список узлов (горизонтальный по “suc”), список альтернатив в список узлов (горизонтальный по “alt”), цикл в узел с указателем “suc” на себя, а “alt” на узел с “empty”, где “alt” = NULL, а “suc” выход из цикла.

5. Построим реализацию в виде структуры данных для рассмотренных в п.3 диаграмм.

6. Программа, реализующая структуру на рис. 1.15 приведена ниже. Булевская переменная b управляет алгоритмом разбора. Переменная state определяет состояние цепочки символов.

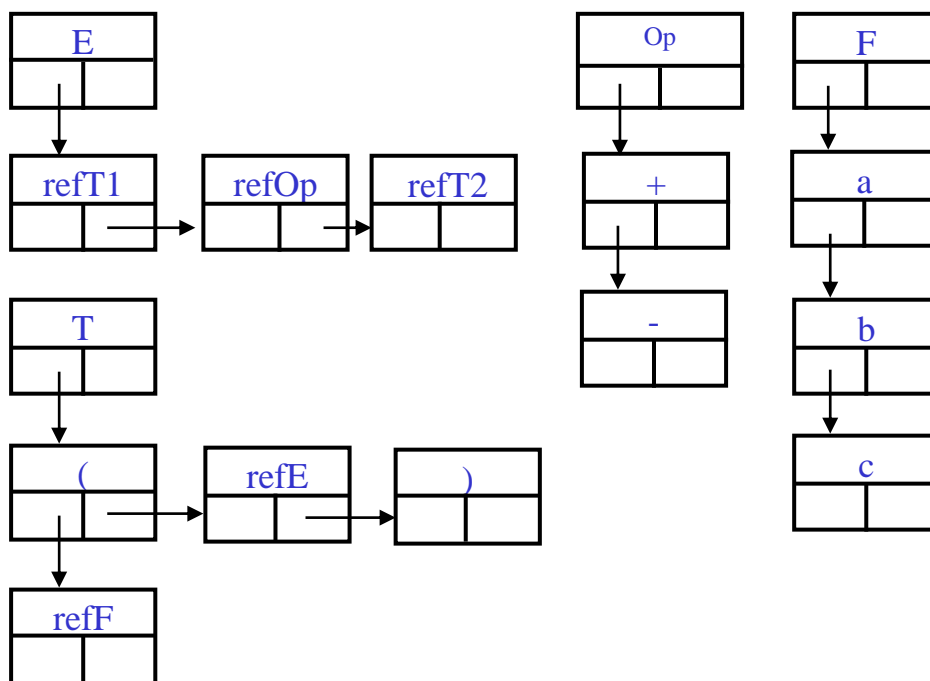


Рис. 1.15. Структура данных для диаграмм п.5

На рис. 1.15. E – начальная вершина Node, а вершины вида refX содержат ссылки на другие вершины. Например, refE содержит ссылку на E, а refT1 и refT2 – на T.

```

using System;

class Program {
    abstract class State { // абстрактный класс с чисто виртуальной функцией
        public abstract void parse(char c);
    }

    // класс для объектов вершин, агрегация по указателю

```

```

class Node {

    public Node (char c){ // задание символа вершинам
        this.alt = null;
        this.suc = null;
        this.sym = null;
        this.c = c;
        this.name = c.ToString();
    }
    public Node(string str) {
        this.alt = null;
        this.suc = null;
        this.sym = null;
        this.c = ' ';
        this.name = str;
    }

    // метод для соединения вершин
    public void link (Node alt, Node suc, Node sym){
        this.alt = alt;
        this.suc = suc;
        this.sym = sym;
    }

    public Node sym; // sym != null
    public Node suc; // terminal empty !=' ' && sym == null
    public Node alt;
    public char c; // empty=' ' && sym == null
    public string name;
}

//подкласс для объекта с состоянием правильного разбора
class OK : State {
    public OK (){}
    public override void parse(char c){
        Console.WriteLine("OK");
    }
}

//подкласс для объекта с состоянием не правильного разбора
class ERROR : State {
    public ERROR(){}
    public override void parse(char c) {
        Console.WriteLine("ERROR string");
    }
}

//класс для автомата агрегация по ссылке объектов классов OK,
// ERROR
class Automate {

    public static State state = null; // полиморфная переменная
    public static ERROR error = new ERROR();
    public static OK ok = new OK();

    public Automate(string str){
        i = 0;
        this.str = str;
    }

    public char getNextChar(){// получить следующий символ из строки
        if (i < str.Length) {
            char ch = str[i];
            i++;
            return ch;
        }
        else {
            return ' ';
        }
    }
}

```

```

        //при окончании строки возвращается пробел
    }

    public virtual void parse(){// начать разбор
        while ( (state !=error)&&(state !=ok) ){
            state.parse(getNextChar()); // принцип подстановки
        }
        state.parse(getNextChar());
    }

    private string str;
    private int i;
} // end class

// определение подкласса для объекта автомата, анализирующего
// контекстно-свободную грамматику
class AutomateCF : Automate {
    public AutomateCF(string str, Node node):base(str) {
        this.node = node; // передача начальной вершины
        c = ' ';
    }

    /// <summary>
    /// Проверяет строку, получаемую из getNextChar(), на соответствие данному
    нетерминалу.
    /// </summary>
    /// <param name="nd">Узел, соответствующий нетерминалу.</param>
    /// <returns></returns>
    public bool parse(Node nd) {
        // Берём первый узел внутри нетерминала.
        Node p = nd.alt;

        // Продолжаем, пока не дойдём до конца нетерминала.
        while (p != null) {
            Console.WriteLine(@"Символ строки = '{0}'. Текущий узел = ""{1}"".",
с, p.name);

            // Переменная содержит true, если текущий узел соответствует
            считанной строке.

            bool b;
            if (p.sym == null) {
                // Текущий символ - не ссылка на нетерминал.

                if (p.c == c) {
                    // Текущий символ совпадает со считанным символом.
                    Console.WriteLine("Совпадение.");

                    b = true;
                    c = getNextChar();
                    state = ok;

                    // конец строки
                    if (c == ' ') {
                        Console.WriteLine("Конец строки.");
                        return true;
                    }
                } else if (p.c == ' ') {
                    // Текущий символ - нетерминал.
                    Console.WriteLine("Начало нетерминала.");
                    b = true;
                } else {
                    // Текущий символ - терминал и не совпадает со
                    считанным символом.

                    Console.WriteLine("Несовпадение.");
                    b = false;
                    state = error;
                }
            }

```

```

    }
    else {
        // Текущий символ - ссылка на нетерминал.
        Console.WriteLine(@"Ссылка на нетерминал = ""{0}""",
p.sym.name);

        // Проверяем, соответствует ли считанная строка текущему
        нетерминалу.
        // Ссылка не может иметь альтернативы, иначе часть
        считанных символов теряется.
        b = parse(p.sym);
    }

    if (b) {
        // Текущий символ соответствует вводу.

        if (p.suc != null) {
            Console.WriteLine(@" Переход в suc = ""{0}""",
p.suc.name);
        }
        else {
            // Текущий символ - последний символ нетерминала.
            Console.WriteLine(" Выход из нетерминала.");
        }

        p = p.suc;
    }
    else {
        // Текущий символ не соответствует вводу.

        if (p.alt != null) {
            Console.WriteLine(@" Переход в alt = ""{0}""",
p.alt.name);
        }
        else {
            // Альтернативы нет -> ввод не соответствует
            нетерминалу.
            Console.WriteLine(" Символ не совпадает с узлом, и
нет альтернативы.");
            return false;
        }

        p = p.alt;
    }
}

// Разбор нетерминала завершился успешно.
return true;
}

public override void parse() { // замещение функции parse в объекте
    c = getNextChar();
    parse(this.node);
    state.parse(getNextChar());
}

private Node node;
private char c;
}

// программа ввода строки с клавиатуры
static string getString () { // различные варианты для тестирования..
    // string str = "c+d-dkf-n";
    // string str = "a+(b-cba)-c";
    // string str = "a+b";
    // string str = "a+b-c-c+a";
    // string str = "a+gfg+fgfg"; // error
    // string str = "a+(b-c)";

```

```

        Console.WriteLine("Введите выражение (например, \"a+(b-c)\"): ");
        string str = Console.ReadLine();
        return str;
    }

    static void Main(string[] args){

        // Построение графа
        // 1. Объявление вершин

        // Вершины - начала нетерминалов
        // E = {T, op, T}
        Node E = new Node("E");
        // T = {(, E, )} | F
        Node T = new Node("T");
        // Op = {+ | -}
        Node Op = new Node("Op");
        // F = {a | b | c}
        Node F = new Node("F");

        // Вершины - ссылки на нетерминалы
        Node refT1 = new Node("refT1");
        Node refOp = new Node("refOp");
        Node refT2 = new Node("refT2");
        Node refE = new Node("refE");
        Node refF = new Node("refF");

        // Вершины - терминалы
        Node braceL = new Node('(');
        Node braceR = new Node(')');
        Node plus = new Node('+');
        Node minus = new Node('-');
        Node c_a = new Node('a');
        Node c_b = new Node('b');
        Node c_c = new Node('c');

        // Соединение вершин в дерево методом link
        E.link(refT1, null, null);
        T.link(braceL, null, null);
        Op.link(plus, null, null);
        F.link(c_a, null, null);

        refT1.link(null, refOp, T);
        refOp.link(null, refT2, Op);
        refT2.link(null, null, T);
        refE.link(null, braceR, E);
        refF.link(null, null, F);

        braceL.link(refF, refE, null);
        braceR.link(null, null, null);
        plus.link(minus, null, null);
        minus.link(null, null, null);
        c_a.link(c_b, null, null);
        c_b.link(c_c, null, null);
        c_c.link(null, null, null);

        //создание объекта автомат
        Automate a = new AutomateCF(getString(), E);
        //инициализация начальных значений
        Automate.state = Automate.ok;
        a.parse(); //синтаксический анализ
    }
}

```

3.6. LL(k)-грамматики

Синтаксический LL-анализатор - анализирует цепочку символов входного алфавита на ленте слева (L) направо, и строит левый (L) вывод грамматики.

Определение 13. КС-грамматика $G = (T, V, P, S)$ без ε -правил называется простой LL(1) грамматикой (s-грамматикой, разделенной грамматикой), если для каждого $v \in V$ все его альтернативы начинаются различными терминальными символами. Единица в названии алгоритма означает, что при чтении анализируемой цепочки, находящейся на входной ленте, входная головка может заглядывать вперед на один символ.

$FIRST(A)$ – это множество первых терминальных символов, которыми начинаются цепочки, выводимые из нетерминала $A \in V$:

$$FIRST(A) = \{a \in T \mid A \Rightarrow_i^+ a\beta, \text{ где } \beta \in (T \cup V)^*\}$$

$FOLLOW(A)$ – это множество следующих терминальных символов, которые могут встретиться непосредственно справа от нетерминала в некоторой сентенциальной форме:

$$FOLLOW(A) = \{a \in T \mid S \Rightarrow_i^* aA\gamma \text{ и } a \in FIRST(\gamma)\}$$

Магазин содержит цепочку $Xa\perp$ (см. рис. 1.17), где Xa – цепочка магазинных символов (X - верхний символ магазина), а символ (\perp) – специальный символ, называемый *маркером дна* магазина. Если верхним символом магазина является *маркер дна*, то магазин пуст. Выходная лента содержит цепочку номеров правил π , представляющую собой текущее состояние левого разбора.

Входная лента

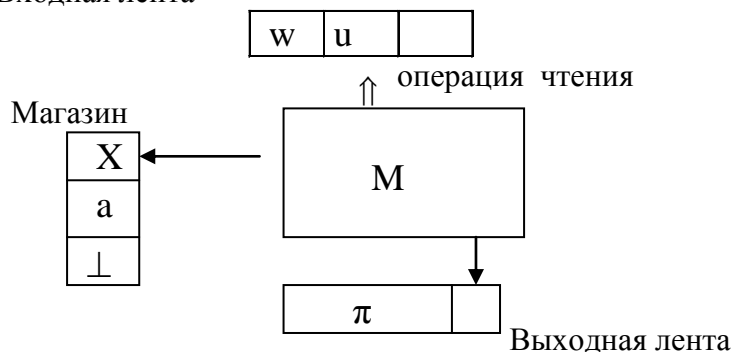


Рис. 1.18. LL(k)-анализатор

Конфигурацию “1-предсказывающего” алгоритма разбора будем представлять в виде $(x, Xa\perp, \pi)$, где x – неиспользуемая часть входной цепочки, Xa – цепочка в магазине, а π – цепочка на выходной ленте, отражающая состояние алгоритма.

Обозначим алфавит магазинных символов (без символа (\perp)) как V_p .

M – управляющая таблица управляет работой алгоритма. M задает отображение множества $(V_p \cup T \cup \{\perp\}) \times (T \cup \{\varepsilon\})$ в множество, состоящее из следующих элементов:

1. (β, i) , где βV_p^* правая часть правила вывода с номером i .
2. ВЫБРОС.
3. ДОПУСК.

4. ОШИБКА.

Работа алгоритма в зависимости от элемента управляющей таблицы $M(X, \alpha) = (\beta, i)$ следующая:

1. $(x, X\alpha, \pi) \vdash (x, \beta\alpha, \pi i)$, если $M(X, \alpha) = (\beta, i)$. В этом случае верхний символ магазина X заменяется на цепочку βV_p^* , и в выходную цепочку дописывается номер правила i . Выходная головка при этом не сдвигается.

2. $(ax, a\alpha, \pi) \vdash (x, \alpha, \pi)$, если $M(a, a) = \text{ВЫБРОС}$. Это означает, что, если верхний символ магазина совпадает с текущим входным символом, он выбрасывается из магазина, и входная головка сдвигается на один символ вправо.

3. Если алгоритм достигает конфигурации $(\varepsilon, \perp, \pi)$, что соответствует элементу управляющей таблицы $M(\perp, \varepsilon) = \text{ДОПУСК}$, то его работа прекращается, и выходная цепочка π является левым разбором входной цепочки.

4. Если алгоритм достигает конфигурации $(x, X\alpha, \pi)$ и $M(X, a) = \text{ОШИБКА}$, то разбор прекращается и выдается сообщение об ошибке.

Конфигурация $(\omega, S\perp, \varepsilon)$, где $S \in V_p$ – начальный символ магазина (начальный символ грамматики), называется *начальной конфигурацией*.

Если $(\omega, S\perp, \varepsilon) \vdash^+ (\varepsilon, \perp, \pi)$, то π называется выходом алгоритма для входа ω .

Алгоритм 3.10. Построение управляющей таблицы M для $LL(1)$ -грамматики

Вход: $LL(1)$ -грамматика $G = (T, V, P, S)$

Выход: Управляющая таблица M для грамматики G .

Таблица M определяется на множестве $(V \cup T \cup \{\perp\}) \times (T \cup \{\varepsilon\})$ по правилам:

1. Если $A \rightarrow \beta$ – правило грамматики с номером i , то $M(A, a) = (\beta, i)$ для всех $a \neq \varepsilon$, принадлежащих множеству $\text{FIRST}(A)$.

Если $\varepsilon \in \text{FIRST}(\beta)$, то $M(A, b) = (\beta, i)$ для всех $b \in \text{FOLLOW}(A)$.

2. $M(a, a) = \text{ВЫБРОС}$ для всех $a \in T$.

3. $M(\perp, \varepsilon) = \text{ДОПУСК}$.

4. В остальных случаях $M(X, a) = \text{ОШИБКА}$ для $X(V \cup T \cup \{\perp\})$ и $a \in T \cup \{\varepsilon\}$

Естественным обобщением $LL(1)$ -грамматик являются $LL(k)$ -грамматики. Для КС-грамматика $G = (T, V, P, S)$ определим множество:

$$\text{FIRST}_k(\alpha) = \{x \mid \alpha \Rightarrow_i^* x\beta \text{ и } |x| = k \text{ или } \alpha \Rightarrow^* x \text{ и } |x| < k\}$$

Применение алгоритма 3.10. требует использование аванцепочек длиной до k символов, что существенно увеличивает размеры управляющей таблицы.

Кроме того, для некоторых $LL(k)$ -грамматик ($k > 1$) верхний символ магазина и аванцепочка длиной k (или меньше) символов не всегда однозначно определяют правило, которое должно быть определено при разборе.

Если в КС-грамматике $G = (T, V, P, S)$ для двух различных A -правил $A \rightarrow \beta$ и $A \rightarrow \gamma$ выполняется:

$$\text{FIRST}_k(\beta \text{ FOLLOW}_k(A)) \cap \text{FIRST}_k(\gamma \text{ FOLLOW}_k(A)) = \emptyset, \text{ то такая КС-грамматика называется сильно } LL(k)\text{-грамматикой.}$$

Сложность построения синтаксических анализаторов и их неэффективность не позволяют практически использовать LL(k)-грамматики.

Проектирование LL(k)-анализатора.

А) Построить управляющую таблицу М для грамматики $G = (T, V, P, S)$, где $P = \{p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8\}$. $p_1: S \rightarrow TE'$, $p_2: E' \rightarrow +TE'$, $p_3: E' \rightarrow \varepsilon$, $p_4: T \rightarrow PT'$, $p_5: T' \rightarrow *PT'$, $p_6: T' \rightarrow \varepsilon$, $p_7: P \rightarrow (S)$, $p_8: P \rightarrow i$

В) Рассмотреть работу алгоритма для выделенной цепочки.

С) Определить является ли LL(k)-грамматика *сильно* LL(k)-грамматикой.

А) Используем алгоритм 3.10. Управляющая таблица должна содержать 11 строк, помеченных символами из множества $(V \cup T \cup \{\perp\})$, и 6 столбцов, помеченных символами из множества $(T \cup \{\varepsilon\})$.

	i	()	+	*	ε
S	TE', 1	TE', 1				
E'			$\varepsilon, 3$	+TE', 2		$\varepsilon, 3$
T	PT', 4	PT', 4				
T'			$\varepsilon, 6$	$\varepsilon, 6$	*PT', 5	$\varepsilon, 6$
P	i, 8	(S), 7				
i	ВЫБРОС					
(ВЫБРОС				
)			ВЫБРОС			
+				ВЫБРОС		
*					ВЫБРОС	
\perp						ДОПУСК

Шаг 1. Строим таблицу построчно. Последовательно рассмотрим все нетерминальные символы.

1. Нетерминалу S соответствует правило вывода грамматики $p_1: S \rightarrow TE'$. Так как $FIRST(TE') = \{ (, i \}$, два терминальных символа, то:

$$M(S, () = M(S, i) = M(TE', 1)$$

2. Для нетерминала E' в грамматике имеются два правила вывода p_2 и p_3 :

$p_2: E' \rightarrow +TE'$ множество $FIRST(+TE') = \{+\}$ и, следовательно, $M(E', +) = M(+TE', 2)$;

$p_3: E' \rightarrow \varepsilon$ имеет простую правую часть, вычислим множество символов, следующих за нетерминалом E' в сентенциальных формах. Построив левый вывод $S \Rightarrow TE' \Rightarrow PT'E' \Rightarrow (S)T'E' \Rightarrow (TE')TE' \dots$, имеем $FOLLOW(E') = \{), \varepsilon \}$.

Таким образом, $M(E',)) = M(E', \varepsilon) = (\varepsilon, 3)$.

Выполняя шаг 1 алгоритма для нетерминалов T, T' и P получим:

Правило грамматики	Множество	Значение М
--------------------	-----------	------------

$p_4: T \rightarrow PT'$ $p_5: T' \rightarrow *PT'$ $p_6: T' \rightarrow \varepsilon$	$FIRST(PT') = \{ (, i \}$ $FIRST(*PT') = \{ * \}$ $FOLLOW(T') = \{ +,), \varepsilon \}$	$M(T', () = M(T', i) = M(TE', 4)$ $M(T', *) = (*PT', 5)$ $M(T', +) = M(T',)) = M(T', \varepsilon) = (\varepsilon, 6)$ $M(P, () = ((S), 7)$ $M(P, i) = (i, 8)$
$p_7: P \rightarrow (S)$ $p_8: P \rightarrow i$	$FIRST((S)) = \{ (\}$ $FIRST(i) = \{ i \}$	

Шаг 2. Далее всем элементам таблицы, находящимся на пересечении строки и столбца, отмеченных одним и тем же терминальным символом, присвоим значение ВЫБРОС.

Шаг 3. Элементу таблицы $M(\perp, \varepsilon)$ присвоим значение ДОПУСК.

Шаг 4. Остальным элементам таблицы присвоим значение ОШИБКА и представим результат в виде таблицы.

Начальное содержимое магазина - $S\perp$

В) Рассмотрим работу алгоритма для цепочки $i+i*i$.

Шаг 1. Алгоритм находится в начальной конфигурации $(i+i*i, S\perp, \varepsilon)$. Значение управляющей таблицы $M(S, i) = M(TE', 1)$, при этом выполняются следующие действия:

- заменить верхний символ магазина S цепочкой TE' ;
- не сдвигать читающую головку;
- на выходную ленту поместить номер использованного правила 1.

Шаг 2. Выполняя действия, аналогичные описанным для шага 1, получим следующие конфигурации:

Текущая конфигурация	Значение М
$(i+i*i, TE'\perp, 1) \mid$	$M(T, i) = (PT', 4)$
$(i+i*i, PT'E'\perp, 14) \mid$	$M(P, i) = (i, 8)$
$(i+i*i, iT'E'\perp, 148) \mid$	$M(i, i) = \text{ВЫБРОС}$

Шаг 3. Алгоритм находится в конфигурации $(i+i*i, iT'E'\perp, 148)$. $M(i, i) = \text{ВЫБРОС}$. Выполняем следующие действия:

- удаляем верхний символ магазина;
- сдвигаем читающую головку на один символ вправо;
- при этом выходная лента не изменяется.

Алгоритм переходит в конфигурацию $(+i*i, T'E'\perp, 148)$.

Выполняя шаги 1 и 2, алгоритм произведет следующие действия:

Текущая конфигурация	Значение М
$(+i*i, T'E'\perp, 148) \mid$	$M(T', +) = (\varepsilon, 6)$
$(+i*i, E'\perp, 1486) \mid$	$M(E', +) = (+TE', 2)$
$(+i*i, +TE'\perp, 14862) \mid$	$M(+, +) = \text{ВЫБРОС}$
$(i*i, TE'\perp, 14862) \mid$	$M(T, i) = (PT', 4)$
$(i*i, PT'E'\perp, 148624) \mid$	$M(P, i) = (i, 8)$
$(i*i, iT'E'\perp, 1486248) \mid$	$M(i, i) = \text{ВЫБРОС}$
$(*i, T'E'\perp, 1486248) \mid$	$M(T', *) = (*PT', 5)$

$(*i, *PT'E'\perp, 14862485) \vdash$	$M(*, *) = \text{ВЫБРОС}$
$(i, PT'E'\perp, 14862485) \vdash$	$M(P, i) = (i, 8)$
$(i, iT'E'\perp, 148624858) \vdash$	$M(i, i) = \text{ВЫБРОС}$
$(\varepsilon, T'E'\perp, 148624858) \vdash$	$M(T', \varepsilon) = (\varepsilon, 6)$
$(\varepsilon, E'\perp, 1486248586) \vdash$	$M(E', \varepsilon) = (\varepsilon, 3)$
$(\varepsilon, \perp, 14862485863)$	

Шаг 5. Алгоритм находится в конфигурации $(\varepsilon, \perp, 14862485863)$

Так как значение $M(\perp, \varepsilon) = \text{ДОПУСК}$, то цепочка $i+i*i$ принадлежит языку и последовательность номеров правил 14862485863 на выходной ленте – это ее разбор.

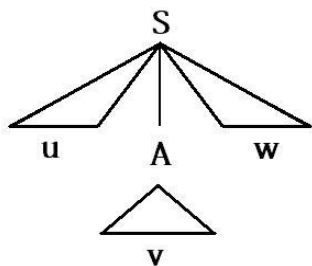
3.7. LR(k)-грамматики

Синтаксический LR-анализатор - анализирует входную цепочку слева направо (L), и строит правый (R) вывод грамматики.

Грамматики, для которых можно построить детерминированный восходящий анализатор, получили название LR(k)-*грамматик* (входная цепочка читается слева (Left) направо, выходом анализатора является правый (Right) разбор, k-число символов входной цепочки, на которое можно “заглянуть” вперед для выделения основы).

Наиболее наглядно LR(k)-*грамматику* можно определить в терминах деревьев вывода.

Определение 14. КС - грамматика $G = (T, V, P, S)$ является LR(k)-*грамматикой*, если просмотрев только часть кроны дерева вывода в этой грамматике, расположенной слева от данной внутренней вершины, часть кроны, выведенную из нее, и следующие k символов входной цепочки, можно установить правило вывода, которое было применено к этой вершине при порождении входной цепочки.



Дерево вывода цепочки uvw.

В определении LR(k)-*грамматики* используется

1. множество $FIRST_k(\gamma)$, состоящее из префиксов длины k терминальных цепочек, выводимых из γ .

Если из γ выводятся терминальные цепочки, длина которых меньше k, то эти цепочки также включаются в множество $FIRST_k(\gamma)$. Формально:

$$FIRST_k(\gamma) = \{x \mid \gamma \Rightarrow^*_1 xw \text{ и } |x| = k \text{ или } \gamma \Rightarrow^*_1 x \text{ и } |x| < k\}$$

2. Пополненной грамматикой, полученной из КС-грамматики $G = (T, V, P, S)$, называется грамматика $G' = (V \cup \{S'\}, T, P \cup \{S' \rightarrow S\}, S')$. Если правила грамматики G' пронумерованы числами 1,2,...,p то, будем считать, что $S' \rightarrow S$ – нулевое правило грамматики G' , а нумерация остальных правил такая же, как в грамматике G .

Определение 15. КС-грамматика $G = (T, V, P, S)$ называется LR(k)-грамматикой для $k \geq 0$, если из существования двух правых выводов для *пополненной* грамматики $G' = (T, V', P', S')$ полученной из G :

$$S' \Rightarrow_r^* \alpha A w \Rightarrow_r \alpha \beta w,$$

$$S' \Rightarrow_r^* \gamma B x \Rightarrow_r \alpha \beta y,$$

для которых $FIRST_k(w) = FIRST_k(y)$ следует, что $\alpha A y = \gamma B x$.

То есть, если $\alpha\beta w$ и $\alpha\beta u$ – правовыводимы цепочки пополненной грамматики G' , у которых $FIRST_k(w) = FIRST_k(u)$ и $A \rightarrow \beta$ – последнее правило, использованное в правом выводе цепочки $\alpha\beta w$, то правило $A \rightarrow \beta$ должно использоваться также в правом разборе при свертке $\alpha\beta u$ к $\alpha A u$.

Поскольку правило грамматики $A \rightarrow \beta$ не зависит от w , то из определения LR(k)-грамматики следует, что во множестве $FIRST_k(w)$ содержится информация достаточная для определения **основы**: кодируемая цепочка символов в верхней части магазина.

Для любой LR(k)-грамматики $G = (T, V, P, S)$ можно построить детерминированный анализатор, который выдает правый разбор входной цепочки.

Анализатор состоит из магазина, входной ленты, выходной ленты и управляющего устройства (пара функций f и g) см. рис. 1.18.

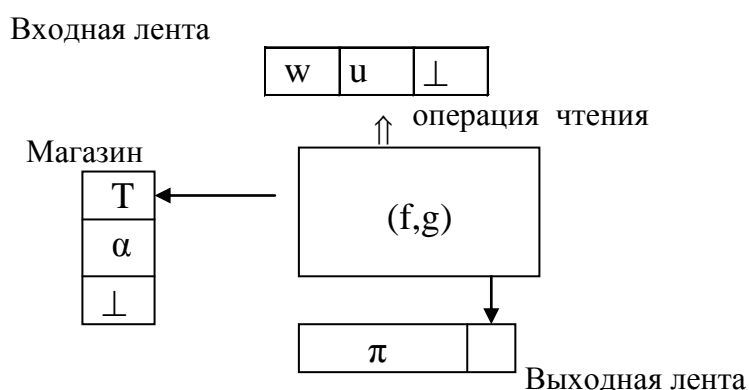


Рис. 1.18. LR(k)-анализатор общий вид

Магазинный алфавит V_p представляет собой множество специальных символов, соответствующих грамматическим вхождениям или их множествам.

Рассмотрим алгоритм “перенос-свертка”, выполняемый детерминированным восходящим синтаксическим анализатором.

Алгоритм состоит в переносе входных символов в магазин до тех пор, пока в его верхней части не встретится основа.

В этот момент производится свертка, в результате которой основа заменяется левой частью основывающего правила. Этот процесс продолжается до тех пор, пока не будет прочитана вся входная цепочка, а в магазине не останется начальный символ грамматики или алгоритм не выдаст сообщение об ошибке.

Пусть в КС-грамматике G имеется правовыводимая цепочка $S \Rightarrow_r^* \alpha x$, такая, что цепочка α заканчивается нетерминальным символом (или это пустая цепочка).

Назовем цепочку α открытой частью цепочки αx , а x – замкнутой частью цепочки. Граница между открытой и замкнутой частью цепочки называется рубежом.

Алгоритм “перенос-свертка” можно рассматривать как программу работы расширенного детерминированного МП-преобразователя.

Рассмотрим правый вывод $S \Rightarrow_r \alpha_0 \Rightarrow_r \alpha_1 \Rightarrow_r \dots \Rightarrow_r \alpha_m = w$. Допустим, что $\alpha_{i-1} = \gamma B z$ и на i -том шаге вывода $\alpha_{i-1} \Rightarrow_r \alpha_i$ было применено правило $B \rightarrow \beta u$

грамматики. Пусть $\gamma V = \alpha A$ и $yz = x$. Тогда цепочка $\alpha_i = \gamma \beta yz = \alpha A x$. При этом преобразователь моделирует обращение вывода $\alpha_{i-1} \leq \alpha_i$.

В начале моделирования i -го шага состояние ДМП-преобразователя определяется цепочкой α_i : в магазине находится αA , а на входной ленте – x . Выполняя разбор, преобразователь будет переносить символы из цепочки x в магазин, пока не определит правый конец основы (возможно, ни одного переноса не потребуется). К этому моменту в магазин будет перенесена цепочка y и в нем будет находиться цепочка $\alpha A y = \gamma B y$.

Преобразователь выполнит свертку, используя правило $B \rightarrow \beta y$, в магазине появится цепочка γV – открытая часть правовыводимой цепочки α_{i-1} , а на входной ленте останется цепочка z – замкнутая часть α_{i-1} .

Моделирование шага вывода завершено.

Два способа построения LR(k) анализаторов

Два способа:

1. Определения активных префиксов - использование расширенного магазинного алфавита (алгоритм перенос-свертка)
2. На основе грамматического вхождения.

Использование расширенного магазинного алфавита: алгоритм перенос-свертка

Входная лента

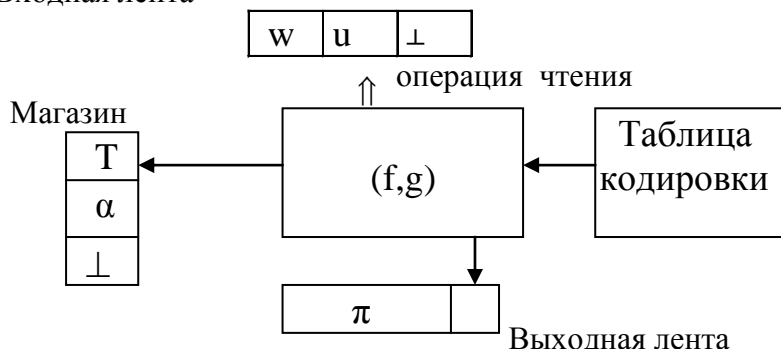


Рис. 1.18. LR(k)-анализатор

Шаг 1. Использование расширенного магазинного алфавита. Анализируя верхний символ магазина определяется наличие в нем основы и правила грамматики для свертки.

1. $S \rightarrow (AS)$
2. $S \rightarrow (b)$
3. $A \rightarrow (SaA)$
4. $A \rightarrow (a)$

Каждый магазинный символ алгоритма содержит сведения о символе грамматики «кодируемой цепочки» (состояние алгоритма). Пример кодирования магазинных символов для рассмотренной ранее грамматики приведен в таблице 3.12.1. Например, магазинный символ « \rangle_3 », находящийся на вершине магазина, означает что:

1. В верхушке магазина находится терминальный символ грамматики « \rangle »;

2. В верхней части магазина находится кодируемая этим символом цепочка символов «(SaA)», то есть *основа* - правая часть правила, 3.

Таблица 3.12.1. Кодирование магазинных символов.

Символ грамматики	Магазинный символ	Кодируемая цепочка
S	S_1	(AS
	S_2	(S
	S_3	$\perp S$
A	A_1	(A
	A_2	(SaA
a	a_1	(Sa
	a_2	(a
b	b_1	(b
($(_1$	(
)	$)_1$	(AS)
	$)_2$	(b)
	$)_3$	(SaA)
	$)_4$	(a)
—	\perp	\perp

Магазинный алфавит построен таким образом, что для каждого магазинного символа (за исключением S_3 и \perp), кодируемая им цепочка является *префиксом правой части* некоторого правила грамматики. И наоборот, каждый непустой префикс правой части правила кодируется некоторым магазинным символом.

Например, правая часть правила грамматики:

$$(1) S \rightarrow (AS)$$

имеет четыре непустых префикса: «(», «(A», «(AS», «(AS)», которые кодируются четырьмя магазинными символами: « $(_1$ », « A_1 », « S_1 », « $)_1$ » соответственно.

Шаг 2. Символ переносится в магазин только в том случае, если он кодирует цепочку, «совместимую» с цепочкой, которая будет находиться в магазине после переноса.

Цепочка, кодируемая данным магазинным символом, *совместима* с цепочкой в магазине, если она является суффиксом магазинной цепочки после переноса данного символа.

Например, если в магазине находится цепочка $\perp(_1S_2a_1$, то алгоритм может втолкнуть в магазин только символ A_2 , так как после этого в магазине будет находиться цепочка $\perp(_1S_2a_1A_2$, суффиксом которой является цепочка, кодируемая символом A_2 , а именно (SaA.

Управление алгоритмом осуществляется при помощи двух функций, приведенных в таблице 8.2, следующим образом:

- Используя значение верхнего символа магазина и входного символа, алгоритм определяет значение функции действия: *ПЕРЕНОС* или *СВЕРТКА*;
- При выполнении переноса определяется значение функции переходов, равное магазинному символу, который нужно втолкнуть в магазин;
- Значение функции действия, равное СВЕРТКА(i), однозначно определяет этот шаг.

Таблица 8.2. Функции управления алгоритмом.

	(функция действий f(u)				функции переходов g(X)					
		a	b)	⊥	S	A	a	b	()
S ₁	П	П	П	П						(₁) ₁
S ₂	П	П	П	П				a ₁		(₁	
S ₃	П	П	П	П	ДОП					(₁	
A ₁	П	П	П	П		S ₁				(₁) ₃
A ₂	П	П	П	П						(₁	
a ₁	П	П	П	П			A ₂		S ₃	(₁	
a ₂	П	П	П	П						(₁) ₄
b ₁	П	П	П	П						(₁) ₂
(₁	П	П	П	П		S ₂	A ₁	a ₂	b ₁	(₁	
) ₁	C(1)	C(1)	C(1)	C(1)	C(1)						
) ₂	C(2)	C(2)	C(2)	C(2)	C(2)						
) ₃	C(3)	C(3)	C(3)	C(3)	C(3)						
) ₄	C(4)	C(4)	C(4)	C(4)	C(4)						
⊥	П	П	П	П		S ₃				(₁	

Пример 3.12.1

Рассмотрим работу алгоритма при разборе цепочки (((b)a(a))((a)(b))). Эта цепочка принадлежит языку, определяемому рассматриваемой грамматикой, так как ее правый вывод:

$$S \Rightarrow_{(1)} (AS) \Rightarrow_{(1)} (A(AS)) \Rightarrow_{(2)} (A(A(b))) \Rightarrow_{(4)} (A((a)(b))) \Rightarrow_{(3)} ((SaA)((a)(b))) \Rightarrow_{(4)} ((Sa(a))((a)(b))) \Rightarrow_{(2)} (((b)a(a))((a)(b)))$$

Шаг 1. Начальная конфигурация алгоритма (\perp , (((b)a(a))((a)(b))), ϵ).

Функция действий f(u) определяется верхним символом магазина и входным символом. Для параметров начальной конфигурации (дно магазина « \perp » и первый входной символ «(» функция действия равна ПЕРЕНОС, f(u)=П.

Определяем магазинный символ, который должен быть помещен в магазин. Так как значение функции переносов для тех же параметров « \perp » и «(» равно магазинному символу «(₁», то алгоритм переходит в конфигурацию:

$$(\perp_1, ((b)a(a))((a)(b))), \epsilon).$$

Шаг 2. Выполняя аналогичные действия алгоритм проходит следующие конфигурации:

$$(\perp, (((b)a(a))((a)(b))), \epsilon) \vdash_s (\perp_1, ((b)a(a))((a)(b))), \epsilon)$$

$$\begin{aligned}
& \vdash_s (\perp_1(\perp_1(b)a(a))((a)(b))), \varepsilon) \\
& \vdash_s (\perp_1(\perp_1(b)a(a))((a)(b))), \varepsilon) \\
& \vdash_s (\perp_1(\perp_1(b_1)a(a))((a)(b))), \varepsilon) \\
& \vdash_s (\perp_1(\perp_1(b_1)_2a(a))((a)(b))), \varepsilon)
\end{aligned}$$

Шаг 3. Алгоритм находится в конфигурации $(\perp_1(\perp_1(b_1)_2a(a))((a)(b))), \varepsilon)$.

Значение функции действий для аргументов « \perp_2 » и символа « a » равно СВЕРТКА(2). Для выполнения свертки по второму правилу грамматики мы должны вычеркнуть из магазина 3 верхних символа и втолкнуть в него некоторый магазинный символ.

Для нахождения значения вталкиваемого символа нужно определить значение функции переходов для следующих аргументов: левой части правила, указанного в значении функции действий и входного символа.

В данном случае это левая часть второго правила грамматики S и входного символа « a », следовательно, получаем значение S_2 , которое помещаем в верхушку магазина.

Шаг 4. Продолжая работу, алгоритм пройдет следующие конфигурации:

$$\begin{aligned}
& (\perp_1(\perp_1(b_1)_2a(a))((a)(b))), \varepsilon) \quad \vdash_{r(2)} (\perp_1(\perp_1S_2a(a))((a)(b))), \perp, 2) \\
& \vdash_s (\perp_1(\perp_1S_2a_1(a))((a)(b))), \perp, 2) \\
& \vdash_s (\perp_1(\perp_1S_2a_1(a))((a)(b))), \perp, 2) \\
& \vdash_s (\perp_1(\perp_1S_2a_1(a_2))((a)(b))), \perp, 2) \\
& \vdash_s (\perp_1(\perp_1S_2a_1(a_2)_4))((a)(b))), \perp, 2) \\
& \vdash_{r(4)} (\perp_1(\perp_1S_2a_1A_2))((a)(b))), \perp, 24) \\
& \vdash_s (\perp_1(\perp_1S_2a_1A_2)_3((a)(b))), \perp, 24) \\
& \vdash_{r(3)} (\perp_1A_1((a)(b))), \perp, 243) \\
& \vdash_s (\perp_1A_1(a)(b))), \perp, 243) \\
& \vdash_s (\perp_1A_1(a)(b))), \perp, 243) \\
& \vdash_s (\perp_1A_1(a_2)(b))), \perp, 243) \\
& \vdash_s (\perp_1A_1(a_2)_4(b))), \perp, 243) \\
& \vdash_{r(4)} (\perp_1A_1(A_1(b))), \perp, 2434) \\
& \vdash_s (\perp_1A_1(A_1(b))), \perp, 2434) \\
& \vdash_s (\perp_1A_1(A_1(b_1))), \perp, 2434) \\
& \vdash_s (\perp_1A_1(A_1(b_1)_2)), \perp, 2434) \\
& \vdash_{r(2)} (\perp_1A_1(A_1S_1)), \perp, 24342) \\
& \vdash_s (\perp_1A_1(A_1S_1)_1), \perp, 24342) \\
& \vdash_{r(1)} (\perp_1A_1S_1), \perp, 243421) \\
& \vdash_s (\perp_1A_1S_1)_1, \perp, 243421) \\
& \vdash_{r(1)} (\perp S_3, \perp, 2434211) \\
& \vdash_s \text{ ДОПУСК.}
\end{aligned}$$

Шаг 5. Алгоритм успешно завершил работу и выдал *правый разбор* входной цепочки $((b)a(a))((a)(b)))$, равный 2434211.

Алгоритм 3.12. LR(k)-алгоритм разбора

Вход: Анализируемая цепочка $z = t_1t_2\dots t_j\dots t_n \in T^*$, где j - номер текущего символа входной цепочки, находящегося под читающей головкой. Управляющая таблица M (множество LR(k)-таблиц) для LR(k)-грамматики $G = (T, V, P, S)$

Выход: Если $z \in L(G)$, то правый разбор цепочки z , в противном случае – сигнал об ошибке.

Алгоритм.

1. $j := 0$.
2. $j := j+1$. Если $j > n$, то выдать сообщение об ошибке и перейти к шагу 5.
3. Определить цепочку u следующим образом:
 - если $k = 0$, то $u = t_j$;
 - если $k \geq 1$ и $j + k - 1 \leq n$, то $u = t_j t_{j+1} \dots t_{j+k-1}$ – первые k -символов цепочки $t_j t_{j+1} \dots t_n$;
 - если $k \geq 1$ и $j + k - 1 > n$, то $u = t_j t_{j+1} \dots t_n$ – остаток входной цепочки.
4. Применить функцию действия f из строки таблицы M , отмеченной верхним символом магазина T , к цепочке u :
 - $f(u) = \text{ПЕРЕНОС (П)}$. Определить функцию перехода $g(t_j)$ из строки таблицы M , отмеченной символом T из верхушки магазина. Если $g(t_j) = T'$ и $T' \in V_p \setminus \{\perp\}$, то записать T' в магазин и перейти к шагу 2. Если $g(t_j) = \text{ОШИБКА}$, то выдать сигнал об ошибке и перейти к шагу 5;
 - $f(u) = (\text{СВЕРТКА}, i) (C)$ и $A \rightarrow \alpha$ – правило вывода с номером i грамматики G . Удалить из верхней части магазина $|\alpha|$ символов, в результате чего в верхушке магазина окажется символ $T' \in V_p \setminus \{\perp\}$, и выдать номер правила i на входную ленту. Определить символ $T'' = g(A)$ из строки таблицы M , отмеченной символом T' , записать его в магазин и перейти к шагу 3.
 - $f(u) = \text{ОШИБКА (О)}$. Выдать сообщение об ошибке и перейти к шагу 5.
 - $f(u) = \text{ДОПУСК (Д)}$. Объявить цепочку, записанную на входной ленте, правым разбором входной цепочки z .
5. Останов.

2. Грамматическое вхождение

Грамматическое вхождение – это символы полного словаря грамматики, снабженные двумя индексами. Первый индекс задает номер i правила грамматики, в правую часть которого входит данный символ, а второй индекс j – номер позиции символа в этой правой части.

Пример: грамматическое вхождение $A_{5,2} = ?$

- (1): $S \rightarrow aAb$
- (2): $S \rightarrow c$
- (3): $A \rightarrow bS$
- (4): $A \rightarrow bB$
- (5): $B \rightarrow aA$

Существуют несколько способов построения детерминированного конечного автомата, распознающего активные префиксы [2, 3, 28, 37]. Рассмотрим вначале достаточно простой способ построения канонической

системы LR(0)-ситуаций, приведенный в [2], для чего определим две функции: CLOSURE и GOTO.

Пример.

Так как LR(0) - ситуация - это просто правило грамматики с точкой в некоторой позиции правой части, то для правила $P \rightarrow (E)$ можно получить 4 ситуации:

$[P \rightarrow \bullet (E)]$

$[P \rightarrow (\bullet E)]$

$[P \rightarrow (E \bullet)]$

$[P \rightarrow (E) \bullet]$

Замечание. LR(0)-ситуация не содержит аванцепочку и, поэтому при ее записи можно опускать квадратные скобки.

Из определения LR(k)-ситуации следует, что для каждого активного префикса существует непустое множество ситуаций. Основная идея построения простого LR(0)-анализатора состоит в том, чтобы вначале построить на базе КС-грамматики детерминированный конечный автомат для распознавания активных префиксов. Для этого ситуации группируются в множества, которые приводят к состояниям анализатора.

Пусть I – множество LR(0)-ситуаций КС-грамматики G . Тогда назовем замыканием множества I множество ситуаций $CLOSURE(I)$, построенное по следующим правилам:

1. Включить в $CLOSURE(I)$ все ситуации из I .
2. Если ситуация $A \rightarrow \alpha \bullet B\beta$ уже включена в $CLOSURE(I)$ и $B \rightarrow \gamma$ - правило грамматики, то добавить в множество $CLOSURE(I)$ ситуацию $B \rightarrow \bullet \gamma$ при условии, что там ее еще нет.
3. Повторять правило 2, до тех пор, пока в $CLOSURE(I)$ нельзя будет включить новую ситуацию.

Второе правило построения можно пояснить следующим образом:

- Наличие ситуации $A \rightarrow \alpha \bullet B\beta$ в множестве $CLOSURE(I)$ говорит о том, что в некоторый момент разбора мы можем встретить в входном потоке анализатора подстроку, выводимую из $B\beta$.
- Если в грамматике имеется правило $B \rightarrow \gamma$, то мы также можем встретить во входном потоке анализатора подстроку, выводимую из γ , следовательно, в $CLOSURE(I)$ нужно включить ситуацию $B \rightarrow \bullet \gamma$.

Рассмотрим расширенную грамматику G'_0 , имеющую следующие правила:

0. $E' \rightarrow E$

1. $E \rightarrow E + T$

2. $E \rightarrow T$

3. $T \rightarrow T * P$

4. $T \rightarrow P$

5. $P \rightarrow i$

6. $P \rightarrow (E)$

Пусть вначале множество ситуации $CLOSURE(I)$ содержит одну ситуацию $E' \rightarrow \cdot E$, т.е. $CLOSURE(I) = \{ E' \rightarrow \cdot E \}$

Ситуация $E' \rightarrow \cdot E$ содержит точку перед символом E , поэтому по второму правилу в $CLOSURE(I)$ необходимо включить E -правила с точкой слева:

$$T \rightarrow \cdot T^* P, T \rightarrow \cdot P, P \rightarrow \cdot i, P \rightarrow \cdot (E)$$

Окончательно получим:

$$CLOSURE(I) = \{ E' \rightarrow \cdot E, E \rightarrow \cdot E + T, E \rightarrow \cdot T, T \rightarrow T^* P, T \rightarrow \cdot P, P \rightarrow \cdot i, P \rightarrow \cdot (E) \}.$$

Если I -множество ситуаций, допустимых для некоторого активного префикса γ , то $GOTO(I, X)$ – это множество ситуаций, допустимых для активного префикса γX .

Аргументами функции $GOTO(I, X)$ являются множество ситуаций I и символ грамматики X .

Определение. Функция $GOTO(I, X)$ определяется как замыкание множества всех ситуаций $[A \rightarrow \alpha X \cdot \beta]$, таких что $[A \rightarrow \alpha \cdot X \beta] \in I$

Пусть для грамматики G_0 множество $I = \{ [E' \rightarrow E \cdot], [E \rightarrow E \cdot + T] \}$.

Для вычисления значения $GOTO(I, +)$ необходимо рассмотреть ситуации, в которых сразу за точкой идет символ $(+)$. Это ситуация $[E \rightarrow E \cdot + T]$. Перемещая точку за символ $(+)$, построим множество ситуаций $\{ [E \rightarrow E + \cdot T] \}$ и замыкание этого множества:

$$\{ T \rightarrow \cdot T^* P, T \rightarrow \cdot P, P \rightarrow \cdot i, P \rightarrow \cdot (E) \}$$

$$\text{Тогда } GOTO(I, +) = \{ [E \rightarrow E \cdot + T], [T \rightarrow \cdot T^* P], [T \rightarrow \cdot P], [P \rightarrow \cdot i], [P \rightarrow \cdot (E)] \}$$

В [2] и [3] приведены алгоритмы. Позволяющие построить каноническую систему множеств $LR(0)$ – ситуаций φ . Процесс построения канонической системы множеств $LR(0)$ – ситуаций можно описать с помощью следующих действий:

$$1. \varphi = \emptyset.$$

2. Включить в φ множество $A_0 = CLOSURE([S' \rightarrow \cdot S])$, которое в начале «не отмечено».

3. Если множество ситуаций A , входящее в систему φ , «не отмечено», то:

- отметить множество A ;
- вычислить для каждого символа $X \in (N \cup \Sigma)$ значение $A' = GOTO(A, X)$;
- если множество $A' \neq \emptyset$ и еще не включено в φ , то включить его в систему множеств как φ «неотмеченное» множество.

4. Повторять шаг 3, пока все множества ситуаций системы φ не будут отмечены.

Пример 8.5.

Определим каноническую систему множеств $LR(0)$ – ситуаций для пополненной грамматики G_0 с правилами:

$$(0) E' \rightarrow E$$

- (1) $E \rightarrow E + T$
- (2) $E \rightarrow T$
- (3) $T \rightarrow T * P$
- (4) $T \rightarrow P$
- (5) $P \rightarrow i$
- (6) $P \rightarrow (E)$

В начале построения система множеств $\varphi = \emptyset$.

Определив множество $A_0 = \text{CLOSURE}(\{|E' \rightarrow \bullet E|\}) = \{E' \rightarrow \bullet E, E \rightarrow \bullet E + T, E \rightarrow \bullet T, T \rightarrow \bullet T * P, T \rightarrow \bullet P, P \rightarrow \bullet i, P \rightarrow \bullet (E)\}$, включим его в систему φ в качестве «неотмеченного» множества:

A_0	$E' \rightarrow \bullet E$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet T * P$ $T \rightarrow \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$
-------	--

Теперь мы должны отметить множество A_0 и определить множества $\text{GOTO}(A_0, X)$ для всех символов грамматики G_0 :

$$A_1 = \text{GOTO}(A_0, E) = \text{GOTO}(\{E' \rightarrow \bullet E, E \rightarrow \bullet E + T\}, E) = \text{CLOSURE}(\{E' \rightarrow E \bullet, E \rightarrow E \bullet + T\}) = \{E' \rightarrow E \bullet, E \rightarrow E \bullet + T\}.$$

Множество A_1 непустое и отсутствует в системе φ . Включим A_1 в φ как «неотмеченное» множество:

A_0	$E' \rightarrow \bullet E$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet T * P$ $T \rightarrow \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$	A_1	$E \rightarrow E \bullet$ $E \rightarrow E \bullet + T$
-------	--	-------	--

Продолжая строить $\text{GOTO}(A_0, X)$, получаем:

$$A_2 = \text{GOTO}(A_0, T) = \text{GOTO}(\{E \rightarrow \bullet T, T \rightarrow \bullet T * P\}, T) = \text{CLOSURE}(\{E \rightarrow T \bullet, T \rightarrow T \bullet * P\}) = \{E \rightarrow T \bullet, T \rightarrow T \bullet * P\}.$$

Включаем A_2 в φ .

$$A_3 = \text{GOTO}(A_0, P) = \text{GOTO}(\{T \rightarrow \bullet P\}, P) = \{T \rightarrow P \bullet\}.$$

Включаем A_3 в φ .

$$A_4 = \text{GOTO}(A_0, () = \text{GOTO}(\{P \rightarrow \bullet (E)\}, () = \{P \rightarrow (\bullet E), E \rightarrow \bullet E + T, E \rightarrow \bullet T, T \rightarrow \bullet T * P, T \rightarrow \bullet P, P \rightarrow \bullet i, P \rightarrow \bullet (E)\}.$$

Включаем A_5 в φ .

Остальные множества $\text{GOTO}(A_0, X)$, где $X \in \{), +, *\}$, пусты, поэтому система φ принимает вид:

A ₀	$E' \rightarrow \bullet E$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet T * P$ $T \rightarrow \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$
A ₁	$E' \rightarrow E \bullet$ $E \rightarrow E \bullet + T$
A ₂	$E \rightarrow T \bullet$ $T \rightarrow T \bullet * P$
A ₃	$T \rightarrow P \bullet$
A ₄	$P \rightarrow (\bullet E)$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet T * P$ $T \rightarrow \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$
A ₅	$P \rightarrow i \bullet$

Продолжая выполнять шаг 3 построения, «отмечаем» множество A₁ и строим множества A₃ = GOTO(A₁, X).

Кроме множества:

$$\begin{aligned}
 A_6 &= \text{GOTO}(A_1, +) = \text{GOTO}(\{E \rightarrow E \bullet + T\}, +) = \\
 &= \text{CLOSURE}(\{E \rightarrow E + \bullet T\}) = \\
 &= \{E \rightarrow E + \bullet T, T \rightarrow \bullet T * P, T \rightarrow \bullet P, P \rightarrow \bullet i, P \rightarrow \bullet (E)\},
 \end{aligned}$$

Остальные множества пусты. Включаем A₆ в φ и получаем:

A ₀	$E' \rightarrow \bullet E$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet T * P$ $T \rightarrow \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$
A ₁	$E' \rightarrow E \bullet$ $E \rightarrow E \bullet + T$

A ₂	E → T • T → T • * P
A ₃	T → P •
A ₄	P → (• E) E → • E + T E → • T T → • T * P T → • P P → • (E) P → • i
A ₅	P → i •
A ₆	E → E + • T T → • T * P T → • P P → • (E) P → • i

Продолжая выполнять шаг 3, получаем:

$$A_7 = \text{GOTO}(A_2, *) = \text{GOTO}(\{E \rightarrow E \cdot * P\}, *) = \text{CLOSURE}(\{T \rightarrow T \cdot * P\} = \{T \rightarrow T \cdot * P, P \rightarrow \cdot i, P \rightarrow \cdot (E)\},$$

Включаем A₇ в φ.

$$A_8 = \text{GOTO}(A_4, E) = \text{GOTO}(\{P \rightarrow (\cdot E), E \rightarrow \cdot E + T\}, E) = \\ = \text{CLOSURE}(\{P \rightarrow (E \cdot), E \rightarrow E \cdot + T\}) = \{P \rightarrow (E \cdot), E \rightarrow E \cdot + T\}.$$

Включаем A₈ в φ.

Вычисляем множество GOTO(A₄, T).

$$A_8 = \text{GOTO}(A_4, T) = \text{GOTO}(\{E \rightarrow \cdot T, T \rightarrow \cdot T * P\}, T) = \\ = \text{CLOSURE}(\{E \rightarrow T \cdot, T \rightarrow T \cdot * P\}) = \{E \rightarrow T \cdot, T \rightarrow T \cdot * P\}.$$

Такое множество (множество A₂) уже есть в φ.

Множества GOTO(A₄, P), GOTO(A₄, () и GOTO(A₄, i) также уже помещены в систему φ.

Множество GOTO(A₅, X) = ∅ для всех X ∈ (N ∪ Σ).

Продолжая аналогичным образом, включаем в систему φ множества:

$$A_9 = \text{GOTO}(A_6, T) = \text{GOTO}(\{E \rightarrow E + \cdot T, T \rightarrow \cdot T * P\}, T) = \\ = \text{CLOSURE}(\{E \rightarrow E + T \cdot, T \rightarrow T \cdot * P\}) = \{E \rightarrow E + T \cdot, T \rightarrow T \cdot * P\},$$

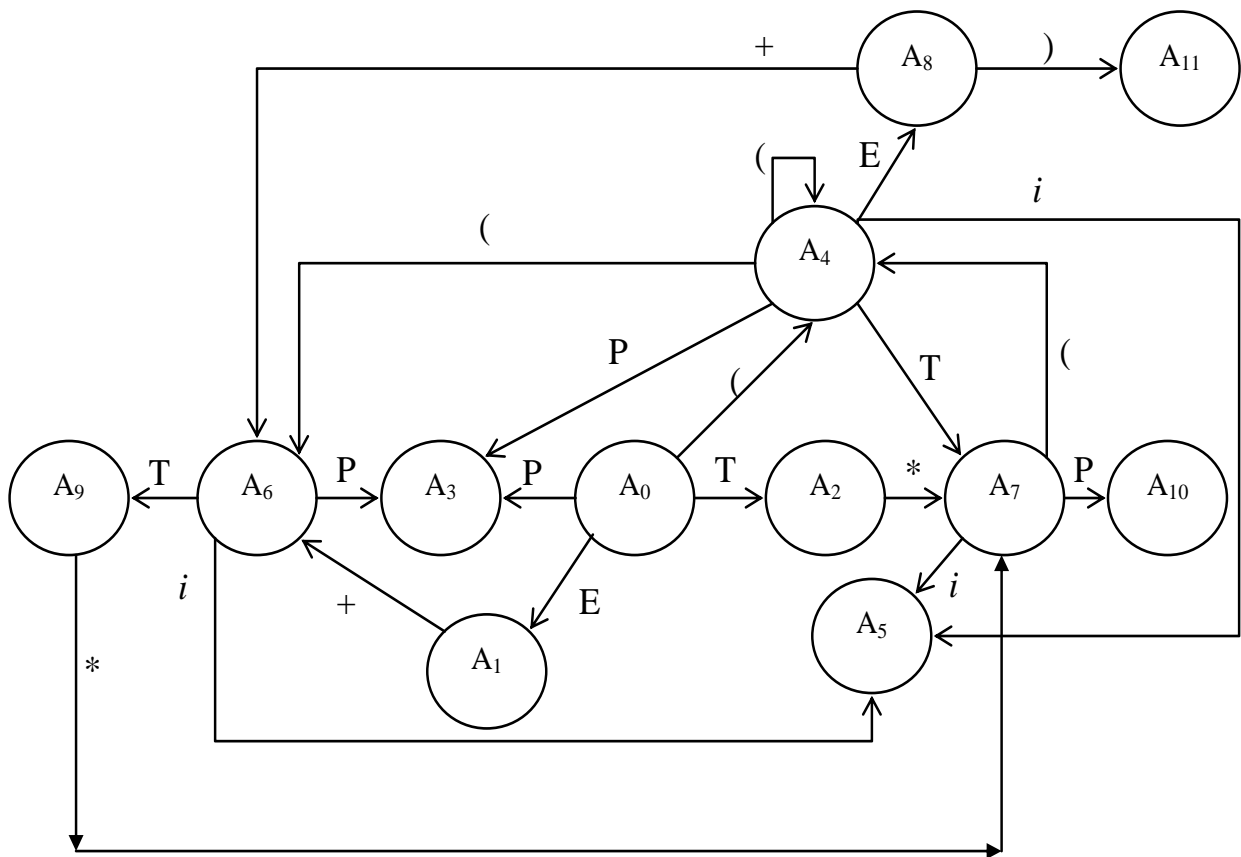
$$A_{10} = \text{GOTO}(A_7, P) = \text{GOTO}(\{T \rightarrow T \cdot * P\}, P) = \\ = \text{CLOSURE}(\{T \rightarrow T \cdot * P \cdot\}) = \{T \rightarrow T \cdot * P \cdot\},$$

$$A_{11} = \text{GOTO}(A_8, () = \text{GOTO}(\{E \rightarrow (\cdot E)\}, () = \\ = \text{CLOSURE}(\{E \rightarrow (E) \cdot\}) = \{E \rightarrow (E) \cdot\}.$$

В результате получаем окончательную систему LR(0) – множеств φ:

A_0	$E' \rightarrow \bullet E$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet T * P$ $T \rightarrow \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$
A_1	$E' \rightarrow E \bullet$ $E \rightarrow E \bullet + T$
A_2	$E \rightarrow T \bullet$ $T \rightarrow T \bullet * P$
A_3	$T \rightarrow P \bullet$
A_4	$P \rightarrow (\bullet E)$ $E \rightarrow \bullet E + T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet T * P$ $T \rightarrow \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$
A_5	$P \rightarrow i \bullet$
A_6	$E \rightarrow E + \bullet T$ $T \rightarrow \bullet T * P$ $T \rightarrow \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$
A_7	$T \rightarrow T * \bullet P$ $P \rightarrow \bullet (E)$ $P \rightarrow \bullet i$
A_8	$E \rightarrow E \bullet + T$ $P \rightarrow (E \bullet)$
A_9	$E \rightarrow E + T \bullet$ $T \rightarrow T \bullet * P$
A_{10}	$T \rightarrow T * P \bullet$
A_{11}	$P \rightarrow (E) \bullet$

Диаграмма переходов ДКА для активных префиксов грамматики G_0 .



Используя каноническую систему $LR(0)$ – множеств, можно представить функцию GOTO в виде диаграммы детерминированного конечного автомата.

Практический интерес представляют $LR(k)$ -грамматики, для которых $k=0$ или $k=1$, в случаях, когда $k > 1$ требуется построение большого числа вспомогательных множеств, поэтому применение $LR(k)$ -анализатора не оправдано.

Алгоритм построения управляющей таблицы M для $LR(0)$ -грамматик основывается на рассмотрении пар грамматических вхождений, которые могут быть представлены соседними магазинными символами в процессе разбора допустимых цепочек.

Пример построения $LR(k)$ -анализатора.

1. Построить управляющую таблицу M для заданной $LR(k)$ грамматики $G = (T, V, P, S)$, где $P = \{p_1, p_2, p_3, p_4, p_5, p_6\}$, $p_1) S \rightarrow aAb$, $p_2) S \rightarrow c$, $p_3) A \rightarrow bS$, $p_4) A \rightarrow Vb$, $p_5) V \rightarrow aA$, $p_6) V \rightarrow c$

2. Описать работу алгоритма.

1. Для построения управляющей таблицы М воспользуемся алгоритмом 1.4.12.

+++++++=

В табл. 3.12.1. приведен результат выполнения алгоритма.

Управляющая таблица М для этой грамматики приведена в табл. 1.3.

Т	f(u) функция действия				g(X)					
	a	b	c	⊥	a	b	c	S	A	B
S ₀				Д						
a ₁	П	П	П		a ₅	b ₃	c ₆		A ₁	B ₄
A ₁	П	П	П			b ₁				
b ₁	C,1	C,1	C,1	C,1						
c ₂	C,2	C,2	C,2	C,2						
b ₃	П	П	П		a ₁		c ₂	S ₃		
S ₃	C,3	C,3	C,3	C,3						
B ₄	П	П	П			b ₄				
b ₄	C,4	C,4	C,4	C,4						
a ₅	П	П	П	П	a ₅	b ₃	c ₆		A ₅	B ₄
A ₅	C,5	C,5	C,5	C,5						
c ₆	C,6	C,6	C,6	C,6						
⊥	П	П	П		a ₁		c ₂	S ₀		

Рассмотрим, например, вторую строку таблицы М, отмеченную грамматическим вхождением a₁, которое представляет собой префикс правой части правила p₁. Пусть b – текущий входной символ. Согласно правилу p₁ за a₁ может следовать грамматическое вхождение нетерминала A₁ или символ, представляющий подцепочку, порождаемую вхождением A₁. Цепочка порождаемая вхождением A₁, может начинаться либо вхождением b₃ (правило p₃), либо вхождением B₄ (правило p₄), которое в свою очередь порождает цепочку, начинающуюся с a₁ или c₆. Следовательно, для входного символа b в магазин можно втолкнуть только грамматическое вхождение b₃, которое является кодированным представлением префикса правой части правила p₃.

Строка таблицы, отмеченная маркером для магазина (⊥), соответствует начальной конфигурации алгоритма. В первый момент времени в магазин можно записать только грамматическое вхождение a₁ или c₂, которые представляют собой префиксы цепочек, выводимых из начального вхождения S₀ (грамматического вхождения начального символа грамматики S, входящего в правую часть нулевого правила вывода пополненной грамматики). Функция действия f(u) не зависит от текущего входного символа, а определяется только верхним символом магазина, то есть для выбора нужного действия (переноса или свертки) алгоритм не должен заглядывать вперед на символы входной цепочки (k=0).

2. Работа алгоритма 3.11. описывается в терминах конфигураций, представляющих собой тройки вида (αТ, ах, π), где αТ – цепочка магазинных символов (Т-верхний символ магазина), ах – необработанная часть входной цепочки, π – выход, построенный к настоящему моменту времени.

Рассмотрим последовательность тактов LR(k)-алгоритма при анализе входной цепочки abcb.

Шаг 1. Начальная конфигурация алгоритма – $(\perp, abcb\perp, \epsilon)$ (магазина находится маркер для магазина, а текущим входным символом является символ a).

Для строки управляющей таблицы, отмеченной символом \perp , $f(a)=\text{ПЕРЕНОС}$ и $g(a) = a_1$, поэтому:

- в магазин записывается символ a_1 (грамматическое вхождение символа a в правую часть первого правила);
- входная головка сдвигается на один символ справа.

Алгоритм приходит в конфигурацию $(\perp a_1, bcb\perp, \epsilon)$.

Шаг 2. Для строки таблицы, отмеченной символом a_1 , $f(b)=\text{ПЕРЕНОС}$ и $g(b) = b_3$, следовательно, алгоритм перейдет в конфигурацию $(\perp a_1 b_3, cb\perp, \epsilon)$. Аналогично для магазинного символа b_3 и текущего символа входной цепочки магазин перейдет в конфигурацию $(\perp a_1 b_3 c_2, b\perp, \epsilon)$

Шаг 3. Рассмотрим строку таблицы M , помеченную грамматическим вхождением c_2 . В этом случае $f(b)=(c, 2)$, значит, необходимо выполнить свертку с использованием правила $p_2) S \rightarrow c$.

Правая часть этого правила содержит только один символ, поэтому:

- удаляем из магазина символ c_2 ;
- определяем значение функции переходов для символа S из левой части правила p_2 в строке таблицы M , отмеченной символом b_3 , который стал верхним символом магазина. Значение функции переходов $g(S) = S_3$.

Алгоритм переходит в конфигурацию $(\perp a_1 b_3 S_3, b\perp, 2)$ и в выходную цепочку записывается число 2 (номер использованного варианта).

Продолжая, получим следующую последовательность тактов:

$$(\perp a_1 b_3 S_3, b\perp, 2) \vdash (\perp a_1 A_1, b\perp, 23) \vdash (\perp a_1 A_1 b_1, \perp, 23) \vdash (\perp a_1 S_0, \perp, 231)$$

Заметим, что конфигурация $(\perp a_1 S_0, \perp, 231)$ является заключительной, а цепочка 231 – правым разбором цепочки $abcb$.

Пусть анализируется цепочка $aabc$, содержащая синтаксическую ошибку, тогда последовательность тактов следующая:

$$\begin{aligned} &(\perp, aabc\perp, \epsilon) \vdash (\perp, aabc\perp, \epsilon) \vdash (\perp a_1, abc\perp, \epsilon) \vdash (\perp a_1 a_5, bc\perp, \epsilon) \\ &\vdash (\perp a_1 a_5 b_3, c\perp, \epsilon) \vdash (\perp a_1 a_5 b_3 c_2, \perp, \epsilon) \vdash (\perp a_1 a_5 b_3 S_3, \perp, 2) \vdash (\perp a_1 a_5 A_5, \perp, 23) \vdash \\ &(\perp a_1 B_4, \perp, 235) \vdash \text{ОШИБКА.} \end{aligned}$$

$f(\epsilon) = \text{ОШИБКА}$. Алгоритм должен выдать сообщение об ошибке и закончит работу.

3.8. Грамматики предшествования

Существует подкласс $LR(k)$ -грамматик, который называется грамматиками предшествования для них легко строится алгоритм типа “перенос-свертка”. При восходящих методах синтаксического анализа в текущей сентенциальной форме выполняется поиск основы α , которая в соответствии с правилом грамматики $A \rightarrow \alpha$ сворачивается к нетерминальному символу A . Основная проблема восходящего синтаксического анализа заключается в поиске основы и нетерминального символа, к которому ее нужно приводить (сворачивать).

Алгоритм 1.4.13. Пусть $G = (T, V, P, S)$ - КС-грамматика, правила которой пронумерованы числами от 1 до p .

Алгоритм типа “перенос-свертка” для грамматики G назовем парой функций $\lambda = (f, g)$, где f – называется функцией переноса, а g – функцией свертки. Функции f и g определяются следующим образом:

1. f отображает $(V \cup T \cup \{\perp\})^* \times (T \cup \{\perp\})^*$ в множество {ПЕРЕНОС, СВЕРТКА, ОШИБКА, ДОПУСК}.

2. g отображает $(V \cup T \cup \{\perp\})^* \times (T \cup \{\perp\})^*$ в множество $\{1, \dots, \text{ОШИБКА}\}$. Если $g(\alpha, w) = i$, то правая часть i -го правила является суффиксом цепочки α .

Пусть имеется сентенциальная форма $\alpha X_1 X_2 \beta$, где $X_1, X_2 \in V \cup T$. На некотором этапе разбора либо символ X_1 , либо символ X_2 , либо они вместе должны войти в основу.

Отношения $(\bullet >)$, $(\bullet =)$ и $(\bullet <)$ называются отношениями предшествования.

Рассмотрим три следующих случая:

X_1 – часть основы, символ X_2 не входит в основу. Символ X_1 будет свернут раньше X_2 . В этом случае говорят, что символ X_1 предшествует X_2 , и записывают как это как $X_1 \bullet > X_2$. X_1 – последний символ основы (хвост основы), т.е. последний символ правой части некоторого правила, а символ X_2 – терминальный символ.

X_1 и X_2 входят в основу, т.е. сворачиваются одновременно, записывают как $X_1 \bullet = X_2$. Символы X_1 и X_2 входят в правую часть некоторого правила грамматики.

X_2 – часть основы, а символ X_1 не входит в основу. Символ X_2 – должен быть свернут раньше X_1 , записывают как $X_1 \bullet < X_2$. X_1 – это первый символ основы (головы основы), т.е. первый символ правой части некоторого правила.

Работу алгоритма типа “перенос-свертка” удобно описывать в терминах конфигураций вида $(\perp X_1 \dots X_m, a_1 \dots a_n, p_1 \dots p_r)$, где :

$\perp X_1 \dots X_m$ – содержимое магазина, X_m – верхний символ магазина и $X_i \in V \cup T$, символ (\perp) для магазина;

$a_1 \dots a_n$ – оставшаяся непрочитанная часть входной цепочки, a_1 – текущий входной символ;

$p_1 \dots p_r$ – разбор, полученный к данному моменту времени.

Один шаг алгоритма λ можно описать с помощью двух отношений: (\vdash^S) (перенос) и (\vdash^r) (свертка), определенных на конфигурациях следующим образом:

1. Если $f(\alpha, aw) = \text{ПЕРЕНОС (П)}$, то входной символ переносится в верхушку магазина и читающая головка сдвигается на один символ вправо. В терминах конфигураций этот процесс описывается так:

$(\alpha, aw, \pi) \vdash^S (\alpha a, w, \pi)$ для $\alpha \in (V \cup T \cup \{\perp\})^*$, $w \in (T \cup \{\varepsilon\})^*$ и $\pi \in \{1, \dots, p\}^*$.

2. Если $f(\alpha\beta, w) = \text{СВЕРТКА (С)}$, $g(\alpha\beta, w) = i$ и $A \rightarrow \beta$ – правило грамматики с номером i , то цепочка β заменяется правой частью правила с номером i , а его номер помещается на выходную ленту, т.е. $(\alpha, aw, \pi) \vdash^r (\alpha A, w, \pi i)$.

3. Если $f(\alpha, w) = \text{ДОПУСК}$, то $(\alpha, w, \pi) \vdash^S \text{ДОПУСК (Д)}$.

В остальных случаях $(\alpha, w, \pi) \vdash^S \text{ОШИБКА}$ (пустое значение в таблице).

Отношение (\vdash) определяется как объединение отношений (\vdash^S) и (\vdash^r) , а транзитивное замыкание отношений (\vdash^+) и (\vdash^*) определяется как обычно.

Для $w \in T^*$ будем записывать $\lambda(w) = \pi$, если $(\perp, w, \varepsilon) \vdash^* (\perp S, \varepsilon, \pi) \vdash^S$ ДОПУСК, и $\lambda(w) = \text{ОШИБКА}$, в противном случае.

Для практического применения алгоритм мало пригоден, так как для определения функции переноса он требует анализа всей цепочки в магазине и всей необработанной части входной цепочки. Аналогичный недостаток имеет алгоритм и при определении функции свертки. В литературе описано несколько вариантов грамматик предшествования.

Для всех классов грамматик предшествования общим является способ поиска правого конца основы (*хвоста* основы). При разборе алгоритмом “перенос-свертка”, всякий раз, когда между верхним символом магазина и первым из необработанных входных символов выполняется отношение $\bullet>$, принимается решение о свертке, в противном случае делается перенос.

Таким образом, с помощью отношения $\bullet>$ локализуется правый конец основы правывыводимой цепочки. Определение левого конца основы (*головы* основы) и нужной свертки осуществляется в зависимости от используемого типа отношения предшествования.

3.9. Пример алгоритм “перенос-свертка”

1. Применить алгоритм типа “перенос-свертка” для заданной грамматики $G = (T, V, P, S)$, где $P = \{p_1, p_2, p_3, p_4\}$, $p_1) S \rightarrow bAb$, $p_2) A \rightarrow cB$, $p_3) A \rightarrow a$, $p_4) B \rightarrow Aad$

2. Описать работу алгоритма.

Функция переноса

f	ax	bx	cx	dx	ε
αA	П	П			
αB	С	С			
Aa	С	С		П	
Ab	П		П		С
Ac	П		П		
Ad	С	С			
\perp		П			
$\perp S$					Д

Функция свертки

g	ax	bx	x	ε
$\perp bAb$				1
αAad		4		
αcB	2	2		
Aa			3	

2. Разберем с помощью построенного алгоритма λ входную цепочку $bcaadb$. Начальная конфигурация алгоритма $(\perp, bcaadb, \varepsilon)$.

Шаг 1. Выполнение алгоритма определяется значением $f(\perp, bcaadb)$, которое, как видно из определения функции f , имеет значение ПЕРЕНОС. Алгоритм переходит в конфигурацию $(\perp b, caadb, \varepsilon)$, выполняя шаг:

$(\perp, bcaadb, \varepsilon) \vdash^r (\perp b, caadb, \varepsilon)$ и следующие шаги:

$(\perp b, caadb, \varepsilon) \vdash^r (\perp bc, aadb, \varepsilon) \vdash^r (\perp bca, adb, \varepsilon)$.

Шаг 2. Очередной шаг выполнения алгоритма определяется значением $f(\perp bca, adb, \varepsilon)$, которое имеет значение СВЕРТКА. Правило, используемое при свертке, определяется значением функции $g(\perp bca, adb) = 3$ и переходит в следующую конфигурацию: $(\perp bca, adb, \varepsilon) \vdash^s (\perp bcA, adb, 3)$, исследующие шаги:

$(\perp bcA, adb, 3) \vdash^r (\perp bcAa, db, 3) \vdash^r (\perp bcAad, b, 3) \vdash^s (\perp bcB, b, 34)$

$\vdash^s (\perp bA, b, 342) \vdash^r (\perp bAb, \varepsilon, 342) \vdash^r (\perp S, \varepsilon, 3421) \vdash^s \text{ДОПУСК}$.

Таким образом, $\lambda(bcaadb) = 3421$, правый вывод цепочки $bcaadb$ должен быть равен 1243. Для проверки построим цепочку по заданному правому выводу:

$S \Rightarrow (p_1) bAb \Rightarrow (p_2) bcBb \Rightarrow (p_4) bcAdb \Rightarrow (p_3) bcaadb$

Глава 4. Грамматики общего вида и машины Тьюринга

4.1. Грамматики общего вида и машина Тьюринга

Определение 15. Грамматика $G = (T, V, P, S)$ называется *грамматикой типа 0*, если на правила вывода не накладывается никаких ограничений (кроме тех, которые указаны в определении грамматики).

Грамматики типа 0 - это самые общие грамматики.

Языки высокого уровня не являются рекурсивными. Проблема выявления принадлежности выражения w языку, порожденному грамматикой типа 0, где $w \in T^*$, в общем случае неразрешима.

Определение 16. Машина Тьюринга (служит для проверки разрешимости алгоритма), в иерархии Хомского занимает самый верхний уровень:

$M = (Q, \Sigma, R, L, B, \delta, q_0)$

Q - множество состояний

Σ - входной алфавит

R, L - правые и левые символы, не входящие в $Q \cup \Sigma$, B - пробел ($B \in \Sigma$)

δ - множество правил **переписывания**

q_0 - исходное состояние

$\delta = q_i a_i \rightarrow q_k a_l, a_i$ - входные символы $q_i a_i \rightarrow q_k \{B, L, R\}$,

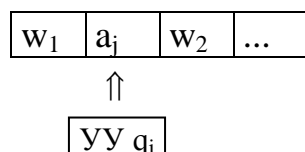
$\delta = Q \times \Sigma \rightarrow Q \times (\Sigma \cup \{B, L, R\})$ $q_i a_j \rightarrow q_k \{B, L, R\}$

Машина Тьюринга - детерминированный преобразователь. Конфигурация машины - множество элементов следующего вида:

$\Sigma^* \times Q \times \Sigma^+$, то есть последовательности $w_1 q_i a_j w_2$, где w_1 и $w_2 \in \Sigma^*$, $a_j \in \Sigma$, $q_i \in Q$.

Правила из δ задают последовательные переходы от одной конфигурации (до тех пор, пока не будет достигнута заключительная конфигурация).

Конфигурация $w_1 q_i a_j w_2$ называется **заключительной**, если к ней неприменимо ни одно из правил δ (т.е. правил с заголовком $q_i a_j$).



Работа машины сводится к следующему, управляющее устройство (УУ):

1. Читает содержимое ячейки, обозреваемой головкой в текущий момент времени.
2. Пишет в обозреваемую ячейку соответствующий символ из входного словаря.
3. Перемещается в соседнюю ячейку, находящуюся слева или справа от обозреваемой ячейки.

Находясь в q_i читает символ a_i и переходит в состояние q_k и либо печатает в этой ячейке символ a_i , либо перемещает головку влево или вправо.

Если к каждой цепочке применить правило, то получится новая конфигурация. Воспринимаемый машиной Тьюринга язык $L = \{w \in \Sigma^* \mid q_0 \text{ и } w \Rightarrow^* C_f, \text{ где } C_f - \text{заключительная конфигурация.}\}$

Если w - цепочка языка, то машина за конечный период времени находит заключительную конфигурацию C_f . Если же $w \notin L$, то машина не останавливается, что приводит к неразрешимости проблемы распознавания принадлежности данного выражения языку.

4.2. Контекстно-зависимые грамматики и ленточные автоматы

Определение 17. Грамматика $G = (T, V, P, S)$ называется *неукорачивающей грамматикой*, если каждое правило из P имеет вид $\alpha \rightarrow \beta$, где $\alpha \subset (T \cup V)^+$, $\beta \subset (T \cup V)^+$ и $|\alpha| \leq |\beta|$.

Определение 18. Грамматика $G = (T, V, P, S)$ называется *контекстно-зависимой (КЗ)*, если каждое правило из P имеет вид $\alpha \rightarrow \beta$, где $\alpha = \xi_1 A \xi_2$; $\beta = \xi_1 \gamma \xi_2$; $A \in V$; $\gamma \subset (T \cup V)^+$; $\xi_1, \xi_2 \subset (T \cup V)^*$.

Грамматику типа 1 можно определить как неукорачивающую либо как контекстно-зависимую.

Выбор определения не влияет на множество языков, порождаемых грамматиками этого класса, поскольку доказано, что множество языков, порождаемых неукорачивающими грамматиками, совпадает с множеством языков, порождаемых КЗ-грамматиками.

4.3. Соотношения между грамматиками и языками

Соотношения между типами грамматик:

1) любая регулярная грамматика является КС-грамматикой;
2) любая регулярная грамматика является укорачивающей КС-грамматикой (УКС). Отметим, что УКС-грамматика, содержащая правила вида $A \rightarrow \varepsilon$, не является КЗ-грамматикой и не является неукорачивающей грамматикой;

3) любая (приведенная) КС-грамматика является КЗ-грамматикой;
4) любая (приведенная) КС-грамматика является неукорачивающей грамматикой;

5) любая КЗ-грамматика является грамматикой типа 0.

6) любая неукорачивающая грамматика является грамматикой типа 0.

Определение 19. Язык $L(G)$ является *языком типа k* , если его можно описать грамматикой типа k .

Соотношения между типами языков:

1) каждый регулярный язык является КС-языком, но существуют КС-языки, которые не являются регулярными (например, $L = \{a^n b^n \mid n > 0\}$);

2) каждый КС-язык является КЗ-языком, но существуют КЗ-языки, которые не являются КС-языками (например, $L = \{a^n b^n c^n \mid n > 0\}$);

3) каждый КЗ-язык является языком типа 0. УКС-язык, содержащий пустую цепочку, не является КЗ-языком. Если язык задан грамматикой типа k , то это не значит, что не существует грамматики типа k' ($k' > k$), описывающей тот же язык. Поэтому, когда говорят о языке типа k , обычно имеют в виду максимально возможный номер k .

Например, КЗ-грамматика $G_1 = (\{0,1\}, \{A,S\}, P_1, S)$ и КС-грамматика $G_2 = (\{0,1\}, \{S\}, P_2, S)$, где $P_1 = \{S \rightarrow 0A1, 0A \rightarrow 00A1, A \rightarrow \varepsilon\}$ и $P_2 = \{S \rightarrow 0S1 \mid 01\}$

описывают один и тот же язык $L = L(G_1) = L(G_2) = \{ 0^n 1^n \mid n > 0 \}$. Язык L называют КС-языком, т.к. существует КС-грамматика, его описывающая. Но он не является регулярным языком, т.к. не существует регулярной грамматики, описывающей этот язык.

Глава 5. Формальные методы описания перевода

- 5.1. Перевод и семантика
- 5.2. СУ-схемы
- 5.3. Транслирующие грамматики
- 5.4. Атрибутные транслирующие грамматики
- 5.5. Методика разработки описания перевода.
- 5.6. Пример разработки АТ грамматики

Глава 6. Разработка и реализация синтаксически управляемого перевода

Рассмотрим два подкласса корректных АТ-грамматик, которые часто используются при проектировании языковых процессоров: *L*-атрибутные транслирующие грамматики и *S*-атрибутные транслирующие грамматики.

1. *L*-атрибутные и *S*-атрибутные транслирующие грамматики

АТ-грамматика называется *L-атрибутной* тогда и только тогда, когда выполняются следующие условия:

- Аргументами правила вычисления значения унаследованного атрибута символа из правой части правила вывода могут быть только унаследованные атрибуты символа из левой части и произвольные атрибуты символов из правой части, расположенные левее рассматриваемого символа.
- Аргументами правила вычисления значения синтезированного атрибута символа из левой части правила вывода являются унаследованные атрибуты этого символа или произвольные атрибуты символов из правой части.
- Аргументами правила вычисления значения синтезированного атрибута операционного символа могут быть только унаследованные атрибуты этого символа.

Является ли произвольная АТ-грамматика *L*-атрибутной, можно проверить, независимо исследуя каждое правило вывода и каждое правило вычисления значения атрибута. Примером *L*-атрибутной транслирующей грамматики является грамматика G^4 .

При выполнении условия 1 определения унаследованные атрибуты каждой вершины дерева вывода зависят (непосредственно или косвенно) только от атрибутов входных символов, расположенных в дереве левее данной вершины, что позволяет использовать *L*-атрибутные грамматики в нисходящих синтаксических анализаторах. Условия 2 и 3 введены с целью сделать АТ-грамматику корректной.

В *L*-атрибутной транслирующей грамматике атрибуты символов *A*, *B* и *C* из правила вывода $A \rightarrow BC$ можно вычислять в следующем порядке:

- унаследованные атрибуты символа *A*;
- унаследованные атрибуты символа *B*;
- синтезированные атрибуты символа *B*;
- унаследованные атрибуты символа *C*;
- синтезированные атрибуты символа *C*;
- синтезированные атрибуты символа *A*.

АТ-грамматика называется *S-атрибутной* тогда и только тогда, когда она является *L*-атрибутной и все атрибуты нетерминалов синтезированные.

Ограничения, накладываемые на *L*-атрибутную транслирующую грамматику, позволяют вычислять значения атрибутов в процессе нисходящего анализа входной цепочки. Нисходящий детерминированный анализатор для *LL(1)*-грамматик требует, чтобы *L*-атрибутная транслирующая грамматика, описывающая перевод, имела форму простого присваивания.

2. Форма простого присваивания

При определении АТ-грамматики в правила вычисления атрибутов записывались в виде операторов присваивания, левые части которых представляют собой атрибут или список атрибутов, а правые части — функцию, использующую в качестве аргументов значения некоторых атрибутов. Простейший случай функции из правой части оператора присваивания — тождественная или константная функция, например $a \leftarrow b$ или $x, y \leftarrow I$.

Правило вычисления атрибутов называется **копирующим правилом** тогда и только тогда, когда левая часть правила — это атрибут или список атрибутов, а правая часть — константа или атрибут. Правая часть называется источником копирующего правила, а каждый атрибут из левой части — приемником копирующего правила.

Если источники нескольких копирующих правил совпадают, то их приемники можно объединить в одну левую часть. Например, правила $z, w \leftarrow a$ и $x, y \leftarrow z$ можно записать в виде: $x, y, z, w \leftarrow a$, т. к. источнику второго правила z , согласно первому правилу, присваивается значение a . Аналогично $x \leftarrow y$ и $a, b \leftarrow y$ можно записать как $a, b, x \leftarrow y$.

Множество копирующих правил называется **независимым**, если источник каждого правила из этого множества не входит ни в одно из других правил множества.

Если два копирующих правила независимы, их нельзя объединять.

Атрибутная транслирующая грамматика имеет форму простого присваивания тогда и только тогда, когда:

- ✓ атрибутов операционных символов.
- ✓ копирующих правил независимо.

Примером АТ-грамматики в форме простого присваивания является грамматика G^4 .

Рассмотрим процедуру преобразования произвольной L-атрибутной транслирующей грамматики в эквивалентную L-атрибутную грамматику в форме простого присваивания. Смысл используемого здесь понятия эквивалентности объясняется далее:

1. Для каждой функции $f(x_1, x_2, \dots, x_n)$, входящей в правило вычисления атрибутов, связанное с некоторым правилом вывода грамматики, создать соответствующий ей операционный символ $\{f\}$, который определяется следующим образом:

$$\begin{array}{l} \{f\} \text{ } x_1 \text{ } x_2, \dots, x_n, p \text{ унаследованные } x_1, x_2, \dots, x_n \\ \text{синтезированный } p \\ p \leftarrow f(x_1, x_2, \dots, x_n) \end{array}$$

2. Для каждого не копирующего правила $z_1, z_2, \dots, z_m \leftarrow f(y_1, y_2, \dots, y_n)$ связанного с некоторым правилом вывода грамматики, найти символы a_1, \dots, a_n , которые не содержатся в этом правиле вывода, и вставить в его правую часть символ $\{f\} a_1, a_2, \dots, a_n, r$. Заменить не копирующее правило на следующие $(n + 1)$ копирующих правил:

$$a_i \leftarrow y_i \text{ для каждого аргумента } y_i,$$

$$z_1, z_2, \dots, z_m \leftarrow r$$

При включении в правило вывода операционного символа необходимо соблюдать следующие условия:

- операционный символ должен располагаться *правее* всех символов правой части правила вывода, атрибутами которых являются аргументы y_1, y_2, \dots, y_n ;
- операционный символ должен располагаться *левее* всех символов правой части правила вывода, атрибутами которых служат z_1, z_2, \dots, z_m ;
- с учетом предыдущих ограничений операционный символ может быть вставлен в любое место правой части правила вывода, но предпочтение следует отдать самой левой из возможных позиций, т. к. это позволяет упростить реализацию синтаксического анализатора.

3. Два копирующих правила, соответствующие одному и тому же правилу вывода, необходимо объединить, если источник одного из них входит в другое. Это достигается удалением правила с лишним источником и объединением его приемников с приемниками оставшегося правила.

Если в качестве источника копирующего правила используется константная функция, являющаяся процедурой-функцией без параметров (например, функция GETNEW из АТ-грамматики), то такие атрибутные правила объединять нельзя, т. к. два разных вызова функции без параметров могут давать разные значения.

Рассмотрим пример преобразования L-атрибутной транслирующей грамматики, порождающей префиксные арифметические выражения над константами, в форму простого присваивания. Атрибутные правила вывода этой грамматики имеют следующий вид:

$$E_p \text{ синтезированный } p$$

$\{\text{ОТВЕТ}\}_p$, унаследованный r

$$(1) \quad S \rightarrow E_p \{\text{ОТВЕТ}\}_r$$

$$r \leftarrow p$$

$$(2) \quad E_p \rightarrow + E_q E_r$$

$$p \leftarrow q + r$$

$$(3) \quad E_p \rightarrow * E_q E_r$$

$$p \leftarrow q * r$$

$$(4) \quad E_p \rightarrow c_r$$

$$p \leftarrow r$$

Входными символами грамматики являются: лексема c , представляющая собой целочисленную константу, и знаки арифметических операций: «+» и «*». Входной символ c имеет один атрибут, значением которого является значение константы. Нетерминальный символ E и операционный символ $\{\text{ОТВЕТ}\}$ также имеют по одному атрибуту. Значением синтезированного атрибута символа E является значение подвыражения, порождаемого этим символом, а значением унаследованного атрибута символа $\{\text{ОТВЕТ}\}$ — значение всего выражения, порождаемого грамматикой.

Исходная грамматика содержит два не копирующих правила: $p \leftarrow q + r$ и $p \leftarrow q * r$, правые части которых представляют собой функции сложения и умножения соответственно. Для преобразования заданной грамматики в форму простого присваивания введем операционные символы $\{\text{СЛОЖИТЬ}\}_{A, B, R}$ и $\{\text{УМНОЖИТЬ}\}_{A, B, R}$ каждый из которых имеет по два унаследованных атрибута A и B и один синтезированный атрибут R . Для операционного символа $\{\text{СЛОЖИТЬ}\}_{A, B, R}$ атрибутивное правило имеет вид: $R \leftarrow A + B$, а для операционного символа $\{\text{УМНОЖИТЬ}\}_{A, B, R}$ — $R \leftarrow A * B$.

Для того чтобы преобразованная грамматика также была L-атрибутной, символ $\{\text{СЛОЖИТЬ}\}$ необходимо поместить правее всех символов правой части правила вывода (2), т. к. одним из аргументов сложения является атрибут самого правого символа E . Атрибут, получающий в качестве своего значения результат сложения, в определении места расположения символа $\{\text{СЛОЖИТЬ}\}$ не участвует, т. к. он не приписан ни к одному из символов правой части. Аналогичные рассуждения относительно операционного символа $\{\text{УМНОЖИТЬ}\}$ определяют крайнюю правую позицию правой части правила (3), как единственно возможное место расположения этого символа. Полученная в результате преобразования L-атрибутная грамматика в форме простого присваивания имеет вид:

E_p синтезированный p

$\{\text{ОТВЕТ}\}_r$, унаследованный r

$\{\text{СЛОЖИТЬ}\}_{A, B, R}$ унаследованный A, B

$R \leftarrow A + B$ синтезированный R

$\{\text{УМНОЖИТЬ}\}_{A, B, R}$ унаследованный A, B

$R \leftarrow A * B$ синтезированный R

$$(1) \quad S \rightarrow E_p \{\text{ОТВЕТ}\}_r$$

$$r \leftarrow p$$

$$(2) \quad E_p \rightarrow + E_q E_r \{\text{СЛОЖИТЬ}\}_{A, B, R}$$

$$A \leftarrow q$$

$$B \leftarrow r$$

$$p \leftarrow R$$

$$(3) \quad E_p \rightarrow * E_q E_r \{\text{УМНОЖИТЬ}\}_{A, B, R}$$

$$A \leftarrow q$$

$$B \leftarrow r$$

$$p \leftarrow R$$

$$(4) \quad E_p \rightarrow c_r$$

$$p \leftarrow r$$

Эта грамматика порождает те же входные цепочки и значение синтезированного атрибута нетерминала E , что и исходная грамматика. Однако формально она не определяет того же самого перевода, т. к. преобразованные правила (2) и (3) *удлиняют активную цепочку*, включив в нее действия, обеспечивающие вычисление функций сложения и умножения соответственно. Для того чтобы преобразованная грамматика определяла тот же перевод, что и исходная, введенные в процессе преобразования операционные символы не следует выдавать в выходную цепочку.

3. Атрибутный перевод для LL(1)-грамматик

Расширим "1-предсказывающий" алгоритм разбора так, чтобы он мог выполнять атрибутный перевод, определяемый L-атрибутной транслирующей грамматикой, входной грамматикой которой является LL(1)-грамматика. Моделирование такого алгоритма можно осуществить с помощью атрибутного ДМП-преобразователя с концевым маркером.

Сначала рассмотрим проблему выполнения синтаксически управляемого перевода, определяемого транслирующей грамматикой цепочечного перевода, входной грамматикой которого является LL(1)-грамматика.

3.1 Реализация синтаксически управляемого перевода для транслирующей грамматики

Преобразуем "1-предсказывающий" алгоритм разбора для LL(1)-грамматик, включив в него действия, обеспечивающие перевод входной цепочки, порождаемой входной грамматикой, в цепочку операционных символов и запись этой цепочки на выходную ленту. Преобразованный таким образом алгоритм в дальнейшем будем называть нисходящим детерминированным процессором с магазинной памятью (нисходящий ДМП-процессор). При этом, если из контекста ясно, что речь идет о нисходящем методе анализа, слово "нисходящий" будем опускать.

В транслирующей грамматике множество терминальных символов разбито на множество входных символов Σ_i , и множество операционных символов Σ_a .

Расширим алфавит магазинных символов, добавив в него операционные символы. Тогда $V_p = \Sigma_i \cup \Sigma_a \cup N$. Затем доопределим управляющую таблицу M , которая для транслирующей грамматики цепочечного перевода задает отображение множества $(V_p \cup \{\perp\}) \times (\Sigma_i \cup \{\epsilon\})$ в множество, состоящее из следующих элементов:

1. (β, y) , где $\beta \in V_p^*$ — цепочка из правой части правила транслирующей грамматики $A \rightarrow y\beta$, а $y \in \Sigma_a^*$;
2. *ВЫДАЧА* (X), где $X \in \Sigma_a$
3. *ВЫБРОС*;
4. *ДОПУСК*;
5. *ОШИБКА*.

Пусть $FIRST(x) = a$, где x — неиспользованная часть входной цепочки. Тогда работу ДМП-процессора в зависимости от элемента управляющей таблицы $M(X, a)$ можно определить следующим образом:

1. $(x, Xa, \pi) \vdash (x, \beta a)$, если $M(X, a) = (\beta, y)$. Верхний символ магазина X заменяется цепочкой $\beta \in V_p^*$, и в выходную цепочку дописывается цепочка операционных символов y . Входная головка при этом не сдвигается.

2. $(x, x\alpha, \pi) \vdash (x, \alpha, \pi X)$, если $M(X, \alpha) = \text{ВЫДАЧА}(X)$. Это означает, что если верхний символ магазина — операционный символ, то он выталкивается из магазина и записывается на выходную ленту. Входная головка не сдвигается.

Действия, соответствующие элементам управляющей таблицы: *ВЫБРОС*, *ДОПУСК* и *ОШИБКА*, остаются теми же, что и в "1-предсказывающем" алгоритме для $LL(1)$ -грамматик.

Опишем алгоритм построения управляющей таблицы M .

Алгоритм построения управляющей таблицы для транслирующей грамматики цепочечного перевода, входной грамматикой которой является $LL(1)$ -грамматика

Вход: Транслирующая грамматика цепочечного перевода $G^T = (N, \Sigma_I, \Sigma_a, P, S)$, входная грамматика которой является $LL(1)$ -грамматика.

Выход: Корректная управляющая таблица M для грамматики G^T .

Описание алгоритма:

■ Управляющая таблица M определяется на множестве $(N \cup \Sigma_I \cup \Sigma_a \cup \{\perp\}) \times (\Sigma_I \cup \{\epsilon\})$ по следующим правилам:

1. Если $A \rightarrow y\beta$ - правило вывода грамматики G^T , где $y \in \Sigma_a^*$, а y — либо ϵ , либо цепочка, начинающаяся с терминала или нетерминала, то $M(A, a) = (\beta, y)$ для всех $a \neq \epsilon$, принадлежащих множеству $\text{FIRST}(\beta)$.

Если $\epsilon \in \text{FIRST}(\beta)$, то $M(A, b) = (\beta, y)$ для всех $b \in \text{FOLLOW}(A)$.

Заметим, что при вычислении $\text{FIRST}(\beta)$ операционные символы, входящие в цепочку β , вычеркиваются.

2. $M(X, a) = \text{ВЫДАЧА}(X)$ для всех $x \in \Sigma_a$ и $a \in \Sigma_i \cup \{\epsilon\}$.

3. $M(a, a) = \text{ВЫБРОС}$ для всех $a \in \Sigma_i$.

4. $M(\perp, \epsilon) = \text{допуск}$

5. В остальных случаях $M(X, a) = \text{ОШИБКА}$ для $x \in (N \cup \Sigma_i \cup \Sigma_a \cup \{\perp\})$ и $a \in \Sigma_i \cup \{\epsilon\}$.

Построим управляющую таблицу для транслирующей грамматики, описывающей перевод инфиксных арифметических выражений в ПОЛИЗ. Эта транслирующая грамматика получена из входной $LL(1)$ -грамматики G , путем включения в нее операционных символов $\{i\}$, $\{+\}$ и $\{*\}$:

- | | |
|--------------------------------|--------------------------------|
| (1) $E \rightarrow TE'$ | (5) $T' \rightarrow *P\{*\}T'$ |
| (2) $E' \rightarrow +T\{+\}E'$ | (6) $T' \rightarrow \epsilon$ |
| (3) $E' \rightarrow \epsilon$ | (7) $P \rightarrow i\{i\}$ |
| (4) $T \rightarrow PT'$ | (8) $P \rightarrow (E)$ |

Управляющая таблица должна содержать 14 строк, помеченных символами из множества $N \cup \Sigma_i \cup \Sigma_a \cup \{\perp\}$, и 6 столбцов, помеченных символами из множества $\Sigma_i \cup \{\epsilon\}$.

Построение управляющей таблицы для строк, отмеченных символами из множества $N \cup \Sigma_i \cup \{\perp\}$, ничем не отличается от построения таблицы для соответствующей входной $LL(1)$ -грамматики (табл. 1.1), а строки управляющей таблицы, отмеченные операционными символами, содержат значения $\text{ВЫДАЧА}(X)$, где $X \in \{i, +, *\}$. Заметим, что элементы таблицы, соответствующие строкам, помеченным операционными символами, для всех столбцов одинаковые, т. к. действия, выполняемые ДМП-процессором в случае, когда верхним символом магазина является операционный символ, не зависят от входного символа.

Таблица 1.1. Управляющая таблица M для транслирующей грамматики, описывающей перевод инфиксных арифметических выражений в ПОЛИЗ

	i	()	+	*	ε
E	TE', ε	TE', ε				
E			ε, ε	+ T{+} E'		ε, ε
T	PT', ε	PT', ε				
T'			ε, ε	ε, ε	*P{*}T', ε	ε, ε
P	i{i}, ε	(E), ε				
1	<i>ВЫБРОС</i>					
(<i>ВЫБРОС</i>				
)			<i>ВЫБРОС</i>			
+				<i>ВЫБРОС</i>		
*					<i>ВЫБРОС</i>	
⊥						<i>ДОПУСК</i>
{i}	<i>ВЫДАЧА ({i})</i>					
{+}	<i>ВЫДАЧА({+})</i>					
{*}	<i>ВЫДАЧА({*})</i>					

Начальное содержимое магазина — $E⊥$

Для входной цепочки $i + i * i$ ДМП-процессор проделает следующую последовательность тактов:

$(i + i * i, E⊥, ε) ⊢ (i + i * i, TE'⊥, ε)$
 $⊢ (i + i * i, PT'E'⊥, ε)$
 $⊢ (i + i * i, \{i\} T'E'⊥, ε)$
 $⊢ (+ i * i, \{i\} T'E'⊥, ε)$
 $⊢ (+ i * i, \{i\} T'E'⊥, \{i\})$
 $⊢ (+ i * i, \{i\} E'⊥, \{i\})$
 $⊢ (+ i * i, +T \{+\} E'⊥, \{i\})$
 $⊢ (i * i, T \{+\} E'⊥, \{i\})$
 $⊢ (i * i, PT' \{+\} E'⊥, \{i\})$
 $⊢ (i * i, i \{i\} T' \{+\} E'⊥, \{i\})$
 $⊢ (* i, \{i\} T' \{+\} E'⊥, \{i\})$
 $⊢ (* i, T' \{+\} E'⊥, \{i\} \{i\})$
 $⊢ (* i, *P\{* \} T' \{+\} E'⊥, \{i\} \{i\})$
 $⊢ (i, P\{* \} T' \{+\} E'⊥, \{i\} \{i\})$
 $⊢ (i, i \{i\} \{* \} T' \{+\} E'⊥, \{i\} \{i\})$
 $⊢ (ε, \{i\} \{* \} T' \{+\} E'⊥, \{i\} \{i\})$
 $⊢ (ε, \{i\} \{* \} T' \{+\} E'⊥, \{i\} \{i\})$
 $⊢ (ε, T' \{+\} E'⊥, \{i\} \{i\} \{i\} \{* \})$
 $⊢ (ε, \{+\} E'⊥, \{i\} \{i\} \{i\} \{* \})$
 $⊢ (ε, E'⊥, \{i\} \{i\} \{i\} \{* \} \{+\})$
 $⊢ (ε, ⊥, \{i\} \{i\} \{i\} \{* \} \{+\})$

ДМП-процессор можно использовать в качестве базового процессора для других видов синтаксически управляемого перевода, если операцию выдачи операционного символа в выходную ленту заменить операциями вызова соответствующих семантических процедур.

3.2 L-атрибутный ДМП-процессор

Процедура преобразования ДМП-процессора в атрибутный ДМП-процессор, которая описывается в данном пособии, требует, чтобы АТ-грамматика, определяющая перевод, имела *форму простого присваивания*.

Рассмотрим построение L-атрибутного ДМП-процессора, реализующего перевод, определяемый L-атрибутной транслирующей грамматикой в форме простого присваивания, входной грамматикой которого является $LL(1)$ -грамматика.

Сначала построим ДМП-процессор, реализующий цепочечный перевод, описываемый транслирующей грамматикой цепочечного перевода, которая получается из заданной L-атрибутной транслирующей грамматики после удаления из нее всех атрибутов. Затем расширим полученный таким образом ДМП-процессор, включив для каждого магазинного символа множество полей для представления атрибутов символа и дополнив управляющую таблицу действиями по вычислению атрибутов и записи их в соответствующие поля.

Для удобства изложения будем считать, что магазинный символ с n атрибутами представляется в магазине $(n + 1)$ -ой ячейками, верхняя из которых содержит имя символа, а остальные — поля для атрибутов. Поля магазинного символа доступны для записи и извлечения атрибутов от момента вталкивания символа в магазин до момента выталкивания его из магазина.

В момент вталкивания символа в магазин в поле для каждого синтезированного атрибута и атрибута входного символа заносится указатель на связанный список полей, соответствующих унаследованным атрибутам, где этот атрибут должен запоминаться, а поле для каждого унаследованного атрибута остается пустым. Содержимое полей синтезированных атрибутов и атрибутов входных символов остается неизменным в течение всего времени нахождения символа в магазине, а поля, соответствующие унаследованным атрибутам, приобретают значения атрибутов к моменту времени, когда магазинный символ окажется в верхушке магазина.

Пусть x — остаток входной цепочки, и $FIRST(x) = a$. Опишем действия, которые должен выполнять L-атрибутный ДМП-процессор в зависимости от элемента управляющей таблицы $M(X, a)$, где X — символ в верхушке магазина.

- **Начальная конфигурация.** В магазине находится маркер дна и начальный символ грамматики. Поля начального символа грамматики, соответствующие унаследованным атрибутам, заполняются начальными значениями атрибутов, которые задаются L-атрибутной транслирующей грамматикой, а в поля, соответствующие синтезированным атрибутам, заносятся пустые указатели, которые служат маркерами конца списков.
- $M(X, a) = \text{ВЫБРОС}$ (символ в верхушке магазина совпадает с текущим входным символом). В этом случае каждое поле верхнего магазинного символа содержит указатель на список полей магазина, в которых требуется поместить значение атрибута текущего входного символа. Операция **ВЫБРОС** расширяется таким образом, что каждый атрибут текущего входного символа копируется во все поля списка на который указывает соответствующее поле верхнего магазинного символа.
- $M(X, a) = \text{ВЫДАЧА}(X)$ (в верхушке магазина находится операционный символ). Операция **ВЫДАЧА}(X)** ДМП-процессора расширяется следующим образом:
 - ✓ унаследованных атрибутов извлекаются из соответствующих полей верхнего магазинного символа и используются затем при выдаче символа в выходную цепочку. При этом следует помнить, что если операционный символ появился в результате преобразования исходной атрибутной транслирующей грамматики в

- форму простого присваивания, то такой символ в выходную цепочку не выдается;
- ✓ значение синтезированных атрибутов вычисляются по правилам вычисления атрибутов, связанным с данным операционным символом, после чего значение каждого синтезированного атрибута помещается во все поля списка, на который указывает соответствующее поле символа из верхушки магазина.
 - $M(X, a) = (\beta, \gamma)$ (в верхушке магазина находится нетерминал). В этом случае L-атрибутный ДМП-процессор вталкивает в магазин цепочку символов β из правой части распознаваемого правила В DSLFTN WTGJXRE операционных символов γ . При этом вычисляются атрибуты операционных символов, которые не вталкиваются в магазин, и заполняются поля атрибутов магазинных символов и символов цепочки β , вталкиваемых в магазин.
 - Источниками атрибутных правил, связанных с правилами вывода L-атрибутной транслирующей грамматики в форме простого присваивания, могут быть только константы, унаследованные атрибуты нетерминалов из левой части правил вывода, атрибуты входных и операционных символов, синтезированные атрибуты нетерминалов из правой части правил вывода. В табл. 1.2 приведены значения источников копирующих правил в момент времени, когда верхним символом магазина является нетерминал.

Таблица 1.2. Значения источников копирующих правил

Источник	Доступ к источнику
Константа	Значение всегда доступно
Унаследованный атрибут нетерминала из левой части правила	Значение источника находится в соответствующем поле верхнего символа магазина
Синтезированный атрибут операционного символа, который не вталкивается в магазин	Значение источника вычисляется в соответствии с атрибутным правилом
Синтезированный атрибут вталкиваемых в магазин нетерминала и операционного символа	Недоступен
Атрибут входного символа	Недоступен

Приемниками атрибутных правил L-атрибутной транслирующей грамматики могут быть только синтезированные атрибуты нетерминалов из левой части правил вывода и унаследо-

ванные атрибуты символов из правой части правил вывода и унаследованные атрибуты символов из правой части правил вывода. В табл. 1.3 приведены поля магазина, соответствующие приемникам атрибутных правил, которые необходимо заполнить во время перехода L-атрибутного ДМП-процессора при $M\{X, a\} = (\beta, y)$.

Таблица 1.3. Значения приемников атрибутных правил

Приемник	Поле
Унаследованный атрибут нетерминала и операционного символа, вталкиваемых в магазин	Соответствующее поле вталкиваемого символа
Синтезированный атрибут нетерминала из левой части правила	Все поля в списке, на который указывает соответствующее поле магазинного символа

При выполнении перехода L-атрибутный ДМП-процессор выполняет следующие атрибутные действия:

1. Вычисляет значения синтезированных атрибутов операционных символов, которые не вталкиваются в магазин, и выдает цепочку операционных символов с их атрибутами в выходную цепочку, если эти символы не появились в результате преобразования грамматики в форму простого присваивания.
2. Если источник копирующего правила доступен, то значение источника помещается в соответствующее поле магазинного символа.
3. Если источник копирующего правила недоступен, то после вталкивания символа в магазин в соответствующие поля символа заносятся указатели на список полей, где будут храниться значения унаследованных атрибутов.

Действия L-атрибутного ДМП-процессора для элементов управляющей таблицы, имеющих значения *ДОПУСК* и *ОШИБКА*, остаются теми же самыми, что у ДМП-процессора.

Построим L-атрибутный ДМП-процессор для L-атрибутной транслирующей грамматики в форме простого присваивания.

На рис. 1.4 показано представление полей магазинных символов, а в табл. 1.5 приведена управляющая таблица для транслирующей грамматики, полученной из исходной L-атрибутной транслирующей грамматики путем вычеркивания из нее атрибутов.

Рисунок 1.4. Разработка и реализация синтаксически управляемого перевода

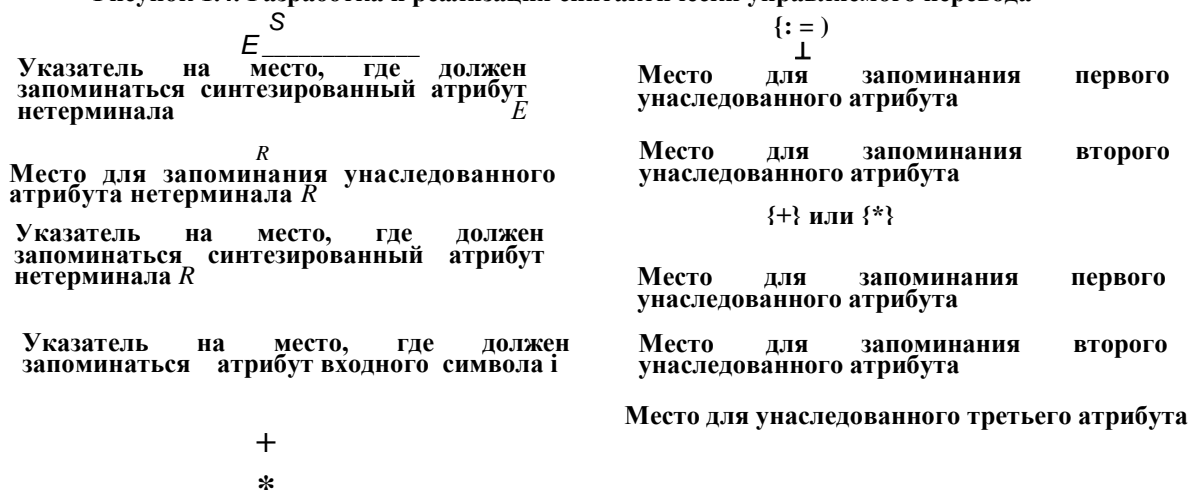


Таблица 1.5. Управляющая таблица для транслирующей грамматики, полученной из исходной грамматики

S	$i := E\{:=\}, \varepsilon$				
E	iR, ε				
R			$+ i\{+\} R, \varepsilon$	$* i\{*\} R, \varepsilon$	ε, ε
/	<i>ВЫБРОС</i>				
:=		<i>ВЫБРОС</i>			
+			<i>ВЫБРОС</i>		
*				<i>ВЫБРОС</i>	
⊥					<i>ДОПУСК</i>
{:=}	<i>ВЫДАЧА($\{:=\}_p q$)</i>				
{+}	<i>ВЫДАЧА($\{+\}_p q, r$)</i>				
{*}	<i>ВЫДАЧА($\{*\}_p q, r$)</i>				

Начальное содержимое магазина — S⊥

4. Атрибутный перевод методом рекурсивного спуска

Сначала рассмотрим, каким образом можно изменить процедуры для распознавания цепочек, порождаемых нетерминальными символами грамматики, для реализации перевода, описываемого транслирующей грамматикой цепочечного перевода.

В этом случае правила составления процедур дополняются следующим правилом: если текущим символом правой части правила вывода грамматики является операционный символ **Y**, ему соответствует вызов процедуры записи операционного символа в выходную цепочку **output (Y)**.

В качестве примера реализации перевода с использованием метода рекурсивного спуска рассмотрим транслирующую грамматику, описывающую перевод инфиксных арифметических выражений в польскую инверсную запись. Эта грамматика построена на основе входной грамматики G_I и имеет следующие правила:

$$\begin{aligned}
 S &\rightarrow E\perp \\
 E &\rightarrow TE' \\
 E' &\rightarrow +T\{+\} E' \mid \varepsilon \\
 T &\rightarrow PT' \\
 T' &\rightarrow *P\{*\} T' \mid \varepsilon \\
 P &\rightarrow (E) \mid i \{i\}
 \end{aligned}$$

Головной модуль **Recurs_Method** и процедуры для распознавания нетерминальных символов **E (Proc_E)** и **T (Proc_T)** не изменятся. Процедура на языке Pascal, реализующая перевод для транслирующих грамматик методом рекурсивного спуска приведена в листинге 1.1.

Листинг 1.1

```

Procedure Recurs_Method_TG (List-Token: tList, List_Oper: tListOper);
                                {List-Token —входная цепочка,
                                List_Oper — выходная цепочка}

Procedure Proc_E
begin

```

```

    Proc_T;
Proc_E1
end {Proc_E};
Procedure Proc_E1;
begin
    if Symb = '+' then
        begin
            NextSymb;
            Proc_T;
            Output
            ({+});
            Proc_E1
        end
    end {Proc_E1}
Procedure Proc_T;
begin
    Proc_P;
    Proc_T1
end
{Proc_T};
Procedure Proc_T1;
begin
    if Symb = '+' then begin
        NextSymb;
        Proc_P;
        Output({*});
        Proc_T1
    end
end {Proc_T1};
Procedure Proc_P;
begin
    if Symb = '('
    then
        begin
            NextSymb;
            Proc_E;
            if Symb = ')' then
                NextSymb;
            else
                Error
            end
        else
            if Symb = 'i'
            then
                begin
                    NextSymb;
                    Output({i})

```


Для того, чтобы методом рекурсивного спуска можно было выполнять L - атрибутный перевод, введем в процедуры распознавания нетерминальных символов *параметр* для каждого атрибута этого символа. При этом если параметру процедуры соответствует унаследованный атрибут нетерминального символа, то вызов этой процедуры производится с фактическим параметром — *значением унаследованного атрибута*, а в случае синтезированного атрибута в качестве фактического параметра используется *переменная*, которой должно быть присвоено *значение синтезированного атрибута в момент выхода из процедуры*.

Для того чтобы процедуры обеспечивали правильную передачу параметров для унаследованных и синтезированных атрибутов, они должны быть написаны на языке программирования, который обеспечивает способы передачи параметров "вызов по значению" (для унаследованных атрибутов) и "вызов по ссылке" (для синтезированных атрибутов).

Замечание. Язык программирования Pascal поддерживает оба метода передачи параметров, а в языке С имеется единственный механизм передачи параметров — "вызов по значению" [36]. Эффект вызова по ссылке в языке С может быть получен, если использовать в качестве параметров указатели.

Поскольку имена атрибутов нетерминальных символов используются в качестве параметров процедур для распознавания этих нетерминалов, при именовании атрибутов необходимо выполнять дополнительное требование, заключающееся в том, что все вхождения некоторого нетерминала в левые части правил вывода АТ-грамматики должны иметь один и тот же список имен атрибутов. Например, в грамматике не может быть таких правил вывода:

$$R_{p1, t2} \rightarrow + i_{q1} \{ * \}_{p2, q2, r1} R_{r2, r1}$$

$$R_{p1, t2} \rightarrow + i_{q1} \{ * \}_{p2, q2, r1} R_{r2, r1}$$

$$R_{p1, p2} \rightarrow \epsilon$$

т. к. первые два вхождения нетерминала R имеют атрибуты p, t , а последнее вхождение — p_1, p_2 . В этом случае необходимо выбрать какой-то один список имен атрибутов нетерминала R (например, p и t), использовать эти имена при описании типа атрибутов этого нетерминала (например, унаследованный p , синтезированный t) и переименовать его атрибуты соответствующим образом.

Указанные ограничения на имена атрибутов не распространяются на атрибуты символов из правых частей правил вывода грамматики.

После того как атрибуты в левых частях правил и в описании их типов переименованы, для упрощения или исключения некоторых правил вычисления атрибутов можно использовать новое соглашение об обозначениях атрибутов, которое формулируется следующим образом: "Если два атрибута получают одно и то же значение, то им можно дать *одно и то же имя* при условии, что для этого не нужно изменять имена атрибутов нетерминала в левой части правила вывода".

Рассмотрим несколько примеров.

$$1. A \rightarrow B_x C_y D_z$$

$$y, z \leftarrow x$$

Атрибутам x, y, z присваивается одно и то же значение, поэтому им можно дать общее имя. Используя новое имя a , получим правило вывода грамматики $(A \rightarrow B_x C_x D_x)$, которое не требует правила вычисления атрибута x .

$$2. A_x \rightarrow B_y f(z) U$$

$$y, z \leftarrow x$$

Атрибутам можно дать одно имя, но оно должно быть x , для того чтобы не изменилось имя атрибута в левой части правила. Выполнив замену имен, получим правило вывода грамматики $A_x \rightarrow B_x \{f\}^*$, при этом отпадает необходимость в атрибутном правиле.

$$3. A_{x,y} \rightarrow a, B_z C_t$$

$$y, z, t \leftarrow x$$

Атрибуты y, z, t, x имеют одно и то же значение, но им нельзя дать одно имя, поскольку нетерминал в левой части правила вывода имеет два атрибута: x и y . В данном случае можно получить лишь частичное упрощение:

$$A_{x,y} \rightarrow a B_y C_y$$

$$y \leftarrow x$$

Новый способ записи имен атрибутов позволяет непосредственно превращать списки атрибутов нетерминальных символов грамматики в списки параметров процедур для распознавания нетерминальных символов. Такой способ записи имеет недостаток, заключающийся в том, что из него не видно, каким образом между атрибутами передается информация. Например, из правила $A \rightarrow B C_x D_x$ не ясно, присваивается ли атрибут нетерминала C атрибуту нетерминала D или наоборот. Если атрибут нетерминала C присваивается атрибуту нетерминала D , то атрибутное правило является L-атрибутным и может использоваться при реализации перевода методом рекурсивного спуска. В противном случае метод рекурсивного спуска неприменим.

Обратившись к описанию типов атрибутов, можно определить порядок передачи информации между атрибутами (значение синтезированного атрибута должно присваиваться унаследованному атрибуту). Однако на практике более удобно использовать обычный способ именования атрибутов и переходить к новому способу записи только после того, как будет доказано, что исходная АТ-грамматика является L-атрибутной.

Для метода рекурсивного спуска не требуется, чтобы АТ-грамматика, описывающая перевод, имела форму простого присваивания.

Опишем детально, как необходимо расширить метод рекурсивного спуска, чтобы он выполнял атрибутный перевод.

Во-первых, изменим процедуру **NextSymbol** таким образом, чтобы она читала очередной символ входной цепочки (лексему) и присваивала класс текущего входного символа переменной **ClassSymb**, а значение лексемы (если оно есть) — переменной **ValSymb**.

Правила составления процедур при условии, что:

- АТ-грамматика, описывающая перевод, является L-атрибутной;
- левые части правил и описания нетерминалов используют одни и те же имена атрибутов;
- для атрибутов, имеющих одно и то же значение, можно использовать одинаковые имена;

дополняются следующими правилами:

- *формальные параметры.* Список имен атрибутов, соответствующий вхождению нетерминала в левые части правил вывода, становится списком формальных параметров соответствующей процедуры;

- *спецификации параметров.* Спецификации атрибутов (УНАСЛЕДОВАННЫЙ или СИНТЕЗИРОВАННЫЙ) переводятся в спецификации формальных параметров по следующим правилам:

- тип УНАСЛЕДОВАННЫЙ соответствует способу передачи параметров "вызов по значению";
- тип СИНТЕЗИРОВАННЫЙ соответствует способу передачи параметров "вызов по ссылке";

- *локальные переменные.* Все имена атрибутов символов данного правила грамматики, кроме тех, что связаны с символом из левой части, становятся локальными переменными соответствующей процедуры;

- *обработка нетерминала из правой части правила.* Для каждого вызова процедуры, соответствующего вхождению нетерминального символа в правую часть правила вывода,

список атрибутов этого вхождения используется в качестве списка фактических параметров;

- *обработка входного символа.* Для каждого вхождения входного символа в правую часть правила вывода грамматики перед вызовом процедуры **NextSymb** в процедуру включается фрагмент кода, который каждой переменной из списка атрибутов входного символа присваивает значение входного атрибута из переменной **ValSymb**;

- *обработка операционного символа.* Для каждого вхождения операционного символа в правую часть правила вывода грамматики в процедуру включается фрагмент кода, который по соответствующим атрибутивным правилам вычисляет значения синтезированных атрибутов операционного символа и присваивает вычисленные значения переменным, соответствующим синтезированным атрибутам. Затем вызывается процедура выдачи операционного символа вместе с атрибутами в выходную строку;

- *обработка правил вычисления атрибутов.* Для каждого правила вычисления атрибутов, сопоставленного правилу вывода грамматики, в процедуру включается фрагмент кода, который вычисляет значение атрибута и присваивает это значение каждой переменной из левой части атрибутивного правила (если соглашение об одинаковых именах выполнено, то в левой части атрибутивного правила будет только один атрибут). Фрагмент кода можно поместить в любом месте процедуры, которое находится:

- после точки, где используемые в правиле атрибуты уже вычислены;
- перед точкой, где впервые используется вычисленное значение атрибута;
- *головной модуль.* Все имена синтезированных атрибутов начального символа грамматики становятся локальными переменными головного модуля. Список фактических параметров вызова процедуры распознавания начального символа грамматики содержит *начальные значения унаследованных атрибутов* и имена синтезированных атрибутов.

Пример 1.

В качестве примера реализации атрибутивного перевода с использованием метода рекурсивного спуска рассмотрим L-атрибутную транслирующую грамматику.

Для повышения наглядности переименуем атрибуты символов грамматики таким образом, чтобы всем вхождениям символов в правые части *разных* правил вывода соответствовали *разные имена атрибутов*, а также выберем одинаковые имена для атрибутов нетерминала **R** из левой части правил вывода с номерами (3), (4) и (5). Пусть *p* — унаследованный атрибут нетерминала **R**, а *t* — его синтезированный атрибут. После преобразования получим следующую грамматику:

E_t синтезированный *t*

$R_{p,t}$ унаследованный *p* синтезированный *t*

Атрибуты операционных символов унаследованные.

(0) $S_0 \rightarrow S \perp$

(1) $S \rightarrow I_a := E_b \{ := \}_{a,b}$

(2) $E_c \rightarrow I_d R_{d,c}$

(3) $R_{p,t} \rightarrow i_{q1} \{ + \}_{p,q1,r1} R_{r1,t}$
 $r_1 \leftarrow \text{GETNEW}$

(4) $R_{p,t} \rightarrow \times i_{q2} \{ * \}_{p,q2,r2} R_{r2,t}$
 $r_2 \leftarrow \text{GETNEW}$

(5) $R_{p,t} \rightarrow \varepsilon$
 $t \leftarrow p$

С учетом выполненных преобразований процедура на языке Pascal, реализующая атрибутивный перевод операторов присваивания некоторого языка программирования в цепочку тетрад с кодами операций: СЛОЖИТЬ, УМНОЖИТЬ, ПРИСВОИТЬ методом рекурсивного спуска, приведена в листинге 2.2.

Листинг 2.2

```
Procedure Recurs_Method_ATG (List_Token: tList, List_Tetr: tListTetr);  
    {List_Token -входная цепочка,  
    List_Tetr - цепочка тетрад}  
  
Procedure Proc_S;  
begin  
    if ClassSymb = ClassId {текущий символ - идентификатор}  
    then  
        begin  
            a := ValSymb;  
            NextSymb;  
            if ClassSymb = ':' then  
                begin  
                    NextSymb;  
                    Proc_E(b);  
                    Output({:=}, a, b)  
                end  
            else  
                Error  
            end  
        end  
    else  
        Error  
    end; {Proc_S}  
  
Procedure Proc_E (var c: integer);  
    var d: integer; {локальная переменная  
    Proc_E}  
    begin  
    if ClassSymb = ClassId {текущий символ - идентификатор}  
    then  
        begin  
            d := ValSymb;  
            NextSymb;  
            Proc_R(d,c)  
        end  
    else  
        Error  
    end;  
{Proc_E}  
  
Procedure Proc_R (p: integer, var t: integer);  
    Var q1, q2, r1, r2: integer; {локальные переменные Proc_R}  
    begin  
        if ClassSymb = '+'  
        then  
            begin  
                NextSymb;  
                if ClassSymb = Classid (текущий символ – идентификатор)  
                then  
                    begin  
                        q2 := ValSymb;
```

```

NextSymb;
GetNew(r2);
Output({+}, p, q2, r2);
Proc_R (r2, t)
end
else
Error;
end
else
t := p
end; {Proc_R}
begin
if ClassSymb = '*'
then begin

```

{ В начале анализа переменная ClassSymb содержит класс первого символа входной цепочки, а переменная ValSymb - значение этого символа}

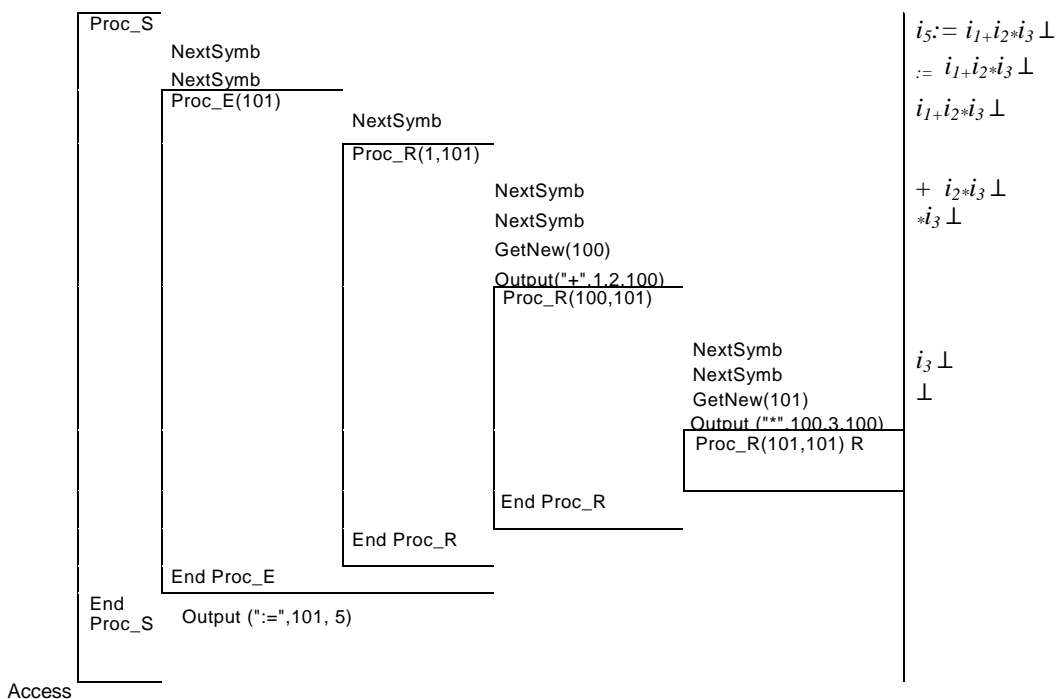
```

Proc_S;
if ClassSymb = '┐'
then
Access
else
Error
end {Recurs_Method_ATG};

```

На рис. 1.3 приведен список имен процедур со значениями фактических параметров в порядке их вызова при переводе входной цепочки $i_5 = i_1 + i_2 * i_3 \perp$ в цепочку тетрад. Справа в строке, соответствующей вызову процедуры **NextSymb**, изображена непрочитанная часть входной цепочки (текущий символ находится слева) непосредственно после вызова этой процедуры.

Рисунок 1.3. Порядок вызова процедур при переводе цепочки $i_5 := i_1 + i_2 * i_3 \perp$



4. S-атрибутный ДМП-процессор

4.1. Математическая модель восходящего ДМП-процессора

Для любого восходящего синтаксического анализатора, рассматриваемого в данном учебнике, последовательность операций переноса и свертки, выполняемых при обработке допустимых входных цепочек, можно описать с помощью транслирующей грамматики. Входной для этой транслирующей грамматики является исходная грамматика, на основе которой построен анализатор. Для получения транслирующей грамматики в самую крайнюю правую позицию каждого i -го правила входной грамматики вставляется операционный символ $\{СВЕРТКА, i\}$. Например, транслирующая грамматика, построенная по входной грамматике $GQ = (\{E, T \mid P\}, \{i, +, *, (,)\}, P, E)$, где $P = \{E \rightarrow E + T, E \rightarrow T, T \rightarrow T * P, P, P \rightarrow i, P \rightarrow (E)\}$, будет выглядеть следующим образом:

- (1) $E \rightarrow E + T \{СВЕРТКА, 1\}$
- (2) $E \rightarrow T \{СВЕРТКА, 2\}$
- (3) $T \rightarrow T * P \{СВЕРТКА, 3\}$
- (4) $T \rightarrow P \{СВЕРТКА, 4\}$
- (5) $P \rightarrow i \{СВЕРТКА, 5\}$
- (6) $P \rightarrow (E) \{СВЕРТКА, 6\}$

Если каждый входной символ в активной цепочке интерпретировать как представление операции **ПЕРЕНОС**, выполняемой в момент времени, когда этот символ является текущим входным символом, то активная цепочка в точности описывает последовательность операций переноса и свертки, выполняемых при обработке входной цепочки. Это объясняется тем, что каждая операция свертки выполняется сразу же после того, как локализована соответствующая *основа* (или первичная фраза), т. е. когда завершается обработка последнего символа в правой части правила вывода. Например, для входной цепочки $i + i * i$, разбор которой рассматривался в разд. 9.4, активная цепочка, порождаемая рассмотренной ранее транслирующей грамматикой, имеет вид:

$i \{СВЕРТКА_6\} + i \{СВЕРТКА_6\} * i \{СВЕРТКА_6\} \{СВЕРТКА_3\} \{СВЕРТКА_1\}$,

что полностью соответствует последовательности операций переноса и свертки, выполняемых при обработке входной цепочки.

Восходящий анализатор можно расширить действиями по выполнению перевода, если перевод определяется постфиксной транслирующей грамматикой. Модификация анализатора заключается в том, что операция свертки расширяется действиями, определяемыми операционными символами соответствующего правила грамматики. Это можно сделать, т. к. при обработке входной цепочки момент времени выполнения свертки для каждого правила грамматики совпадает с моментом выполнения действий по переводу для этого правила. Например, для постфиксной транслирующей грамматики цепочечного перевода:

- $$\begin{aligned} E &\rightarrow E + T \{+\} \\ E &\rightarrow T \\ T &\rightarrow T * P \{*\} \\ T &\rightarrow P \\ P &\rightarrow i \{i\} \\ P &\rightarrow (E) \end{aligned}$$

операции свертки расширяются следующим образом: **СВЕРТКА_1** будет обеспечивать выдачу операционного символа $\{+\}$ в выходную строку, **СВЕРТКА_3** — выдачу операционного символа $\{*\}$, а **СВЕРТКА_6** — выдачу операционного символа $\{i\}$.

Синтаксический анализатор, дополненный формальными действиями по выполнению перевода, принято называть восходящим ДМП-процессором.

4.2 Реализация S-атрибутного ДМП-процессора

Восходящий ДМП-процессор можно легко преобразовать в S-атрибутный ДМП-процессор реализующий атрибутный перевод, определяемый постфиксной S-атрибутной транслирующей грамматикой.

В S-атрибутном ДМП-процессоре каждый магазинный символ имеет конечное множество полей для представления атрибутов. Так же, как и в L-атрибутном ДМП-процессоре, примем, что магазинный символ с атрибутами представляется в магазине $(n + 1)$ -ой ячейками, верхняя из которых содержит имя символа, остальные — поля для атрибутов. Поля магазинного символа, предназначенные для атрибутов, *заполняются* значениями атрибутов в момент *вталкивания символа в магазин* и не изменяются до момента выталкивания его из магазина.

В S-атрибутном ДМП-процессоре *операция переноса* расширяется таким образом, что значения атрибутов переносимого входного символа помещаются в соответствующие поля вталкиваемого при переносе магазинного символа.

При выполнении *операции свертки* для правила с номером i верхние символы магазина представляют собой правую часть i -го правила вывода входной грамматики, а поля магазинных символов содержат значения атрибутов соответствующих символов грамматики.

Расширенная операция свертки использует эти значения для вычисления значений всех атрибутов операционных символов, связанных с правилом вывода транслирующей грамматики, и значений всех атрибутов нетерминала из левой части правила. *Значения атрибутов операционных символов* используются для выдачи результатов в выходную ленту или выполнения других действий, определяемых этими символами. *Атрибуты нетерминала из левой части правила вывода* записываются в соответствующие поля магазинного символа, который соответствует этому нетерминалу и вталкивается в магазин во время свертки.

Приложение А. Порядок выполнения лабораторных работ

Практические занятия по курсу состоят из тематических работ, включающих одну или несколько лабораторных работ.

Лабораторная работа №1-3

Спроектировать автоматную грамматику по заданному языку L , построить конечный автомат.

1. Изучить классификацию Хомского (см. раздел 1.1., 1.2., 1.3.). Ответьте на вопрос, какие грамматики называются автоматными. Какие есть виды автоматных грамматик.

2. Спроектировать по заданному языку L автоматную грамматику и конечный автомат. Используйте пример, и последовательность выполнения работы из раздела 2.3. “Практическая работа 1”.

1. Постановка задачи.
2. Входные и выходные данные.
3. Спроектировать грамматику (Лаб 1).
4. Определить свойства грамматики.

5. Спроектировать конечный автомат, составить диаграмму переходов КА и реализовать на C# (Лаб 2.).

5.1. Создать проект – консольное приложение:

5.2. Используйте следующие фрагменты программ для реализации КА, распознающего заданный язык.

Пример 1. Создайте проекта простого консольного приложения.

```
using System;
using System.Collections.Generic;
using System.Collections;
using System.Linq;
using System.Text;

namespace ConsoleApplication {
    class Program {
        static void Main(string[] args) {
            ArrayList Q = new ArrayList();
            string str;
            Console.WriteLine("Hello");
            str = Console.ReadLine();
            Console.WriteLine("is= "+str);
            Console.ReadLine();
        }
    }
}
```

```
}
```

Пример 2. Разбор входной строки по символам.

Пример 3. Проектирование и реализация правил.

```
    Delta delta = null;
// 1. transition
    delta = new Delta("S0", "0", new ArrayList() { "A","qf" });
    DeltaList.Add(delta);
// 2. transition
    delta = new Delta("A", "1", new ArrayList() { "B" });
    DeltaList.Add(delta);
// 3. transition
    delta = new Delta("B", "0", new ArrayList() { "A","qf" });
    DeltaList.Add(delta);

class Delta { // структура Delta правил переписывания
    private string LeftNoTerm = null;
    private string LeftTerm    = null;
    private ArrayList Right    = null;
    public string leftNoTerm { get { return LeftNoTerm; } set { LeftNoTerm =
value; } }
    public string leftTerm { get { return LeftTerm; } set { LeftTerm = value; } }
    public ArrayList right { get { return Right; } set { Right = value; } }

    // модель правила
    // delta( A,          1)          = {qf}
    //      LeftNoTerm LeftTerm      Right
    public Delta(string LeftNoTerm, string LeftTerm, ArrayList Right) {
        this.LeftNoTerm = LeftNoTerm;
        this.LeftTerm = LeftTerm;
        this.Right = Right;
    }

    public void DebugDeltaRight() {
        int i = 0;
        for (; i < this.right.Count; i++) {
            if (i == 0)
                System.Console.Write("(" + this.right[i]);
            else
                System.Console.Write(", " + this.right[i]);
        }
        if (i == 0)
```

```

        System.Console.WriteLine();
    else
        System.Console.WriteLine(" ");
    }
} // end class Delta

```

Пример 4. Определение автомата как объекта и тестирование.

```

class Automate { // NDKA (Q,Sigma,deltaList,q0,F)
    ArrayList Q      = null;    // множество всех состояний
    ArrayList Sigma  = null;    // конечный алфавит входных символов
    ArrayList DeltaList = new ArrayList(); // множество всех правил
    string Q0       = null;    // начальное состояние
    ArrayList F      = null;    // заключительные состояния
    // атрибутное программирование на C#
    public ArrayList q { get { return Q; } set { Q = value; } }
    public ArrayList sigma { get { return Sigma; } set { Sigma = value; } }
    public ArrayList deltaList { get { return DeltaList; }
                                set { DeltaList = value; } }
    public string q0 { get { return Q0; } set { Q0 = value; } }
    public ArrayList f { get { return F; } set { F = value; } }

    public Automate() {
        this.Q      = new ArrayList();
        this.Sigma  = new ArrayList();
        this.DeltaList = new ArrayList();
        this.F      = new ArrayList();
    }

    public Automate(string aname) {
        System.Console.WriteLine(aname);
        // сделать диалог, инициализирующий NDKA
        // альтернативные состояния переходов хранить в массиве см. Test
        //init();
        Test();
    }

    public void Test() { // задание правил для тестирования
        Q = new ArrayList() { "S0", "A", "B", "qf" }; // "C" для отладки
        Sigma = new ArrayList() { "0", "1" };
        q0 = "S0";
        F = new ArrayList() { "qf" };
        ... задать правила
    } // end test

```


Пример 5. Считывание символа и выбор правила:

```
foreach (Delta d in this.DeltaList) {  
    if (d.leftNoTerm == q && d.leftTerm == chain.Substring(i,1) ) {  
        ...  
    }  
    ...  
}
```

6. Определить свойства КА. Реализовать алгоритм преобразования НДКА в ДКА (Лаб 3.).

Пример 6. Построение побитового кода булеана для множества мощности n :

Для построения Булеана вначале строим все побитовые комбинации от 0 до $2^n - 1$, где n – мощность множества. Затем побитовая комбинация преобразуется в символы. Например, бинарная последовательность для множества для множества из символов {SAB}, $n = 3$ приведена в таблице. "1" означает, что символ в подмножестве, "0" – его отсутствие:

Пример 7. Построение по побитовому коду дельта правила:

Реализация алгоритма преобразования на C#:

```
// NDKA (Q,S,Delta,q0,F)  
namespace NDKA2DKA {  
    class Program {  
        static void Main(string[] args) {  
            Automate NDKA = new Automate("NDKA");  
            NDKA.Debug();  
            Converter converter = new Converter();  
            Automate DKA = converter.convert(NDKA);  
            DKA.Debug();  
            DKA.recognize(DKA.q0, "01010", 0);  
        }  
    }  
  
    //  
    class Converter {  
        Automate DKA = null;  
        // множество всех правил deltaListAll  
        ArrayList deltaListAll = null; // all transitions  
        // подмножества, которые содержат все заключительные  
        // состояния qf, то есть F'  
        ArrayList FAll = null;  
  
        public Converter() {}  
    }  
}
```

```

public Automate convert(Automate NDKA) {
    // инициализировать данные для каждого вызова convert
    this.DKA = new Automate();
    this.deltaListAll = new ArrayList();
    this.FAll = new ArrayList();

    // Шаг 1. Init q0 & sigma
    DKA.q0 = NDKA.q0;
    DKA.sigma = NDKA.sigma;

    //2. Создать множество всех подмножеств по Q (булеан) и
    //3. Создать множество всех правил DeltaList
    // 4 и 5, подмножества, которые содержат заключительные
    // состояния qf, то есть F'
    BuildDeltaList(NDKA);
    // Шаг 6. Определить достижимые состояния,
    // исключить недостижимые состояния из множества Q'
    // 1. Берем начальное состояние и определяем правило дельта,
    Reachability(DKA.q0);
    return DKA;
}

void BuildDeltaList(Automate NDKA) {
    ArrayList right = null; // для нового правила
    int count = NDKA.q.Count;
    // Time Complexity: O(n^2*n), Space Complexity: O(1)
    // 1. set size of power set of a set with set size n is (2**n)
    int sizeOfPowerSet = (int)Math.Pow(2, count);
    string leftNoTerm = null; // is subset
    string[] noTerm = null; // для split
    // 2. Run from counter 000..0 to 111..1
    Console.WriteLine("Boolean_____");
    for (int counter = 0; counter < sizeOfPowerSet; counter++) {
        leftNoTerm = null;
        for (int j = 0; j < count; j++) {
            // Check if j-th bit in the counter is set If set then build set, use comma
            // Console.WriteLine("! 0x{0:x8}", 1 << j);
            if ((counter & 1 << j) != 0) {
                // System.Console.WriteLine("NDKA.q[j] = " + NDKA.q[j]);
                if (leftNoTerm != null) leftNoTerm = leftNoTerm + ',' + NDKA.q[j];
                else leftNoTerm = "" + NDKA.q[j];
            }
        }
        if (leftNoTerm != null) { // 2**n -1 без пустых подмножеств
            // Найти delta'(S, a)

```

```

        noTerm = leftNoTerm.Split(',');
        // Шаг 4. построить subset F
        BuildFAll(leftNoTerm, NDKA.f);
        foreach (string leftTerm in NDKA.sigma){
            // по deltaListNDKA посмотреть имеющиеся правила для данного, одно
            foreach (string n in noTerm) { // ищем правило
                right = findTransition(n, leftTerm, NDKA.deltaList);
                if (right != null) {
                    deltaListAll.Add(new Delta(leftNoTerm, leftTerm, new ArrayList(right)));
                    break;
                }
            }
        }
        System.Console.WriteLine("      "+leftNoTerm); // булеан
    }
    } // end for
    DebugDeltaList(deltaListAll);
    DebugF(FAll);
    } // BuildDeltaList

// найти переход
public ArrayList findTransition (string leftNoTerm, string leftTerm,
                                ArrayList NDKAdeltaList) {
    foreach (Delta d in NDKAdeltaList) { // найдено правило в ArrayList
        if (d.leftNoTerm == leftNoTerm && d.leftTerm == leftTerm)
            return d.right;
    }
    return null;
}

void BuildFAll(string leftNoTerm, ArrayList qf) {
    // если в подмножестве noTerm есть заключительное состояние NDKA.f
    string[] noTerm = leftNoTerm.Split(',');
    foreach (string n in noTerm) { // ищем правило
        foreach (string f in qf) {
            if (n == f) {
                FAll.Add(leftNoTerm);
                // System.Console.WriteLine(" FAll.Add = " + leftNoTerm);
                return;
            }
        }
    }
} // end BuildFAll

void Reachability(string q) {

```

```

// 1. Берем состояние по правилу дельта и определяем следующее дельта
string right = null;
foreach (Delta d in deltaListAll) {
    if (d.leftNoTerm == q) {
        // преобразовать в метку подмножество из right
        d.right = markSubset(d.right);
        DKA.deltaList.Add(d);
        DKA.q.Add(q);
        // всегда один элемент, так как markSubset
        right = d.right[0].ToString();
        break;
    }
}
if (right == null) return; // нет достижимых состояний
if (DKA.q.Contains(right)) { // это состояние уже было, останов
    // заключительное состояние должно быть F'
    if (Fall.Contains(right))
        // в F' оставить последнее достижимое состояние
        DKA.f.Add(right);
    else {
        System.Console.WriteLine(" Reachability error " + 01);
        return;
    }
}
else Reachability(right);
} // end Reachability

ArrayList markSubset(ArrayList right) {
    string r = null;
    foreach (string s in right) {
        if (r != null) r = r + ',' + s;
        else r = s;
    }
    return new ArrayList(){r};
}

void DebugF(ArrayList F) {
    System.Console.WriteLine(" F all:_ ");
    for (int i = 0; i < F.Count; i++) {
        System.Console.WriteLine("      "+F[i]);
    }
}

public void DebugDeltaList(ArrayList deltaList) {
    System.Console.WriteLine("deltaList all:_ ");
}

```

```

foreach (Delta d in deltaList) {
    Console.WriteLine("      (" + d.leftNoTerm + ")," + d.leftTerm + " = ");
    d.DebugDeltaRight();
}
}
} // end class Convertor

```

7. Оформить работу согласно указанным шагам.

Лабораторная работа №4-6

Привести заданную КС-грамматику $G = (T, V, P, S)$ к приведенной форме.

1. Изучить алгоритмы приведения КС-грамматик к приведенной форме. (см. раздел 3.4.). Ответьте на вопрос, какие КС-грамматики называются грамматиками в приведенной форме. Лаб. 4 А,В. Лаб. 5 С,Д., Лаб. 6 F,Е.

2. Использовать алгоритмы преобразования КС-грамматик..

А). Устранить из грамматики G бесполезные символы. Применить алгоритм 3.3. к грамматике G .

В). Устранить из грамматики G ϵ -правила, применить алгоритм 3.5.

С). Устранить из КС грамматики G цепные правила, применить алгоритм 3.6.

Д). Устранить левую рекурсию в заданной КС-грамматике G_1 , порождающей скобочные арифметические выражения. Применить алгоритм 3.7. к грамматике G .

Ф). Определить в какой форме (Грейбах, Хомского) находится КС-грамматика G' .

Е). G' – приведенная КС-грамматика.

3. Оформить работу согласно шагам.

Лабораторная работа №7-8

Построить МП-автомат P и расширенный МП-автомат по КС-грамматике $G = (T, V, P, S)$, без левой рекурсии. Написать последовательность тактов автоматов для выделенной цепочки. Определить свойства автоматов. Лаб. 7. -2.1. Лаб. 8. -2.2.

1. Изучить алгоритмы построения МП-автомат P и расширенного МП-автомата по заданной КС-грамматике (см. раздел 3.4.). Ответьте на вопрос, чем отличается МП-автомат P от расширенного МП-автомата.

2. Выполнить построение согласно алгоритмам. Смотрите пример и последовательность выполнения работы из раздела 3.4. Практическая работа 4.

2.1. А). Построить МП-автомат по КС-грамматике G , используя алгоритм 3.8. Для моделирования магазина используйте Stack из библиотеки C#.

```
Stack <string> stack = new Stack <string> ();
```

```
stack.Push("c");  
stack.Push("d");
```

```
Console.WriteLine(" simbol: " + stack.Pop());
```

В). Определить последовательность тактов МП-автомата для выделенной цепочки.

2.2. А). Построить расширенный МП-автомат, используя алгоритм 3.9.

В). Определить последовательность тактов расширенного МП-автомата при анализе входной выделенной цепочки Р.

3. *Определить свойства построенный МП-автоматов.*

4. *Оформить работу согласно шагам.*

Лабораторная работа №9-10

Разработать контекстно-свободную грамматику по заданной строке (см. раздел 3.3.). Алгоритм разбора реализуется в виде процедур, каждой, из которой соответствует диаграмма.

1. *Повторить классификацию Хомского* (см. раздел 1.1.). Ответьте на вопрос, какие грамматики называются контекстно-свободными. Какие есть виды контекстно-свободных грамматик.

2. *Спроектировать по заданной строке* контекстно-свободную грамматику и автомат с магазинной памятью. Используйте пример и последовательность выполнения работы из раздела 3.3. Обратите внимание, что в качестве неявного стека могут выступать: рекурсивная процедура и древовидная структура объектов. Лаб. 9 1-4. Лаб. 10 5-7.

1. Спроектировать контекстно-свободную грамматику.

2. Записать вывод заданной строки по грамматики.

3. Определить свойства грамматики.

4. Устранить левую рекурсию и записать вывод заданной строки по грамматики.

5. Оптимизировать грамматику метасимволами {...} и записать вывод заданной строки по грамматики.

6. Составить синтаксический граф для грамматики.

7. Преобразовать граф в программу.

```
#include "stdafx.h"  
#include <iostream>  
using namespace std;
```

```
class State { // абстрактный класс с чисто виртуальной функцией  
public:  
    virtual void parse(char c)=0;  
};  
//подкласс для объекта с состоянием правильного разбора  
class OK:public State {
```

```

    public:
        OK () {}
        virtual void parse(char c){cout << "OK" << endl;}
};
//подкласс для объекта с состоянием не правильного разбора
class ERROR:public State {
    public:
        ERROR() {}
        virtual void parse(char c){cout << "ERROR string"<< endl;}
};
//класс для автомата агрегация по ссылке объектов классов OK,
ERROR
class Automate {
    public:
        static State* state; // полиморфная переменная
        static ERROR* error;
        static OK* ok;

        Automate(string String){
            i=0;
            this->String = String;
        }
        char getNextChar(){// получить следующий символ из строки
            if (i<(int)String.length()){char ch=String[i];i++;return ch;}
            else {return ' ';}
            //при окончании строки возвращается пробел
        }
        virtual void parse(){ // начать разбор
            while ( (state !=error)&&(state !=ok) ){
                state->parse(getNextChar());
            }
            state->parse(getNextChar()); // принцип подстановки
        }
    private:
        string String;
        int i;
}; // end class
// присвоение начальных значений переменным автомата
State* Automate::state = NULL;
ERROR* Automate::error = new ERROR;
OK* Automate::ok = new OK;

// определение подкласса для объекта автомата, анализирующего
// контекстно-свободную грамматику
class AutomateCF: public Automate {
    public:
        AutomateCF(string String):Automate(String) {c = ' ';}
        virtual void parse(){// замещение функции parse в объекте
            // класса Automate
            c = getNextChar(); // вызов функции класса автомат
            E(); // вызов метода
            state->parse(getNextChar());
        }
}

```

```

void T () { // реализация функции по диаграмме
    cout << "step1 T=" << c << endl;
    if ( (c == 'a' || c == 'b' || c == 'c')) {
        c = getNextChar();
        cout << "step2 T=" << c << endl;
    }
    if (c == ')') {c=getNextChar();}
    else if (c == '(') {c=getNextChar(); E();}
}

void E () { // реализация функции по диаграмме
    cout << "step0 E=" << c << endl;
    T(); // вызов функции
    cout << "step1 E=" << c << endl;
    while ( (c == '+' || c == '-') ) {
        c = getNextChar(); T(); // вызов функции
        cout << "step2 E=" << c << endl;
    }
    if (c == ' ') {state=ok;}
    else {state=error;}
}
private:
    char c;
};

// программа ввода строки с клавиатуры
string getString () {
    //string String = "c+d-dkf-n";
    string String = "a+(b-c)";
    // string String = "a+b";
    char c;
    int N = (int)String.length();
    cout << "Enter string " << N << " char " << endl;
    for (int i=0; i < N; i++) {
        cin >> c;
        String[i] = c;
    }
    return String;
}

int main() {
    //создание объекта автомат
    Automate * a = new AutomateCF(getString());
    //инициализация начальных значений
    a->state = Automate::ok;
    a->parse(); //синтаксический анализ
    return 0;
}

```

Обратите внимание, что в переменную “с”, записывается следующий символ, который доступен из функций T() и E(). В функциях указаны трассировочные шаги выполнения, которые необходимо сохранить в программе. Состояния OK и ERROR автомата используются для идентификации выполнения разбора.

3. Оформить работу согласно шагам.

Лабораторная работа №11-10

Реализовать контекстно-свободную грамматику, полученную в работе 5 на основе, таблично-управляемой программы грамматического разбора. Изучить контекстно-зависимые грамматики (см. раздел 1.4.).

1. Повторить классификацию Хомского (см. раздел 1.1.). Ответьте на вопрос, какие грамматики называются контекстно – зависимыми, неограниченными. Какие свойства контекстно – зависимых и неограниченных грамматик, в чем их отличие. Чем они отличаются от изученных ранее грамматик ?

2. Разработать объектно-ориентированную реализацию для таблично-управляемой программы грамматического разбора.

1. Представить граф грамматики в виде структуры данных (см. раздел 3.3.).

2. Классифицировать типы вершин.

3. Выполнить подстановку графов и получить как можно меньшее число графов.

4. Последовательность вершин графа преобразовать в структуру узлов.

5. Реализовать программу разбора по структуре узлов.

В алгоритме разбора, результат каждого шага разбора выводится на экран, чтобы можно было видеть, как происходит разбор.

```
#include "stdafx.h"
#using <mscorlib.dll>

using namespace System;
using namespace std;
class State { // абстрактный класс с чисто виртуальной функцией
public:
    virtual void parse(char c)=0;
};
// класс для объектов вершин, агрегация по указателю
class Node {
public:
    Node (char c){ // задание символа вершинам
        this->alt = NULL;
        this->suc = NULL;
        this->sym = NULL;
        this->c = c;
    }
    // метод для соединения вершин
    void link (Node *alt, Node *suc, Node *sym){
        this->alt = alt;
        this->suc = suc;
        this->sym = sym;
    }
};
```

```

    }

    Node * sym; // sym != NULL
    Node * suc; // terminal empty !=' ' && sym == NULL
    Node * alt;
    char c;      // empty=' ' && sym == NULL
};

//подкласс для объекта с состоянием правильного разбора
class OK:public State {
public:
    OK () {}
    virtual void parse(char c){cout << "OK" << endl;}
};
//подкласс для объекта с состоянием не правильного разбора
class ERROR:public State {
public:
    ERROR () {}
    virtual void parse(char c){cout << "ERROR string"<< endl;}
};

//класс для автомата агрегация по ссылке объектов классов OK,
// ERROR
class Automate {
public:
    static State* state; // полиморфная переменная
    static ERROR* error;
    static OK* ok;

    Automate(string String){
        i=0;
        this->String = String;
    }
    char getNextChar(){// получить следующий символ из строки
if (i<(int)String.length()){char ch=String[i]; i++; return ch;}
    else {return ' ';}
    //при окончании строки возвращается пробел
    }
    virtual void parse(){// начать разбор
        while ( (state !=error)&&(state !=ok) ){
            state->parse(getNextChar()); // принцип подстановки
        }
        state->parse(getNextChar());
    }
private:
    string String;
    int i;
}; // end class

// присвоение начальных значений переменным автомата
State* Automate::state = NULL;
ERROR* Automate::error = new ERROR;
OK* Automate::ok = new OK;

```

```

// определение подкласса для объекта автомата, анализирующего
// контекстно-свободную грамматику
class AutomateCF: public Automate {
public:
    AutomateCF(string String, Node &node):Automate(String) {
        this->node = &node; // передача начальной вершины
        c = ' ';
    }

    void parse(Node * nd, bool b){
        Node *p = nd->alt; // Node head
        do {
            if (p->sym == NULL){ // терминал или пусто "empty"
                if (p->c == c){
                    cout<< "step01 c=" << c<< ": p->c ="<<p->c<<endl;
                    b=true; c=getNextChar();
                    state=ok;
                    // конец строки
                }
                if (c == ' ') {cout<< "step012 c="<<endl; return;}
                } else if (p->c == ' ') { // "empty"
                    cout<< "step02 if c=" << c<< ": p->c ="<<p->c<<endl;
                    b=true;
                } else { b=false; state=error;
                    cout<< "step03 c=" << c<< ": p->c ="<<p->c<<endl;
                }
            }
            else if (p->sym != NULL){ //- nil
                cout<< "step04 c=" << c<< "p->sym ="<<p->sym<<endl;
                parse(p->sym, b);
            }
        }
        if (b) {
            cout<< "step05 ok: next c=" << c<< endl;
            p=p->suc;
        }
        else {p=p->alt;}
    } while (p != NULL);
} // end parse

virtual void parse(){ // замещение функции parse в объекте
    c = getNextChar();
    parse(this->node, true);
    state->parse(getNextChar());
}

private:
    Node * node;
    char c;
};

// программа ввода строки с клавиатуры
string getString (){ // различные варианты для тестирования..
    // string String = "c+d-dkf-n";
    // string String = "a+(b-cba)-c";

```

```

// string  String = "a+b";
// string  String = "a+b-c-c+a";
// string  String = "a+gfg+fgfg"; // error
    string  String = "a+(b-c)";
char c;
int N = (int)String.length();
cout << "Enter string "<<N<<" char "<< endl;
for (int i=0; i < N;i++){
    cin>>c;
    String[i] = c;
}
return String;
}

int main(){
    // Построение графа
    // 1. Объявление вершин
    Node E = Node(' ');
    Node braceL = Node('(');
    Node nil = Node(' ');
    Node braceR = Node(')');
    Node c_a = Node('a');
    Node c_b = Node('b');
    Node c_c = Node('c');
    Node emptyR = Node(' ');
    Node plus = Node('+');
    Node minus = Node('-');
    Node nil_1 = Node(' ');

    // 2. Соединение вершин в дерево методом link
    // void link (Node * alt, Node * suc, Node * sym)
    E.link(&braceL,NULL,NULL);
    braceL.link(&c_a,&nil,NULL);
    nil.link(NULL,&braceR,&E);
    braceR.link(NULL,&plus,NULL);
    c_a.link(&c_b,NULL,NULL);
    c_b.link(&c_c,NULL,NULL);
    c_c.link(&emptyR,NULL,NULL);
    emptyR.link(NULL,&plus,NULL);
    plus.link(&minus,&nil_1,NULL);
    minus.link(NULL,&nil_1,NULL);
    nil_1.link(NULL,NULL,&E);

    //создание объекта автомат
    Automate * a = new AutomateCF(getString(),E);
    //инициализация начальных значений
    a->state = Automate::ok;
    a->parse(); //синтаксический анализ

    return 0;
}

```

6. Оформить работу согласно шагам.

Лабораторная работа №12-13.

Построить управляющую таблицу М для LL(k)-грамматики, написать правило вывода, определить является ли G грамматика *сильно* LL(k)-грамматикой (см. раздел 3.4.).

1. Изучить раздел 3.4. Построение управляющей таблицы М для грамматики $G = (T, V, P, S)$, работу алгоритма для определенной цепочки и определение *сильно* LL(k)-грамматики.

Пример построения управляющей таблицы.

```
using System;
using System.Collections.Generic;
using System.Collections;
using System.Linq;
using System.Text;

namespace lab {
    class Program {
        static void Main(string[] args) {
            MTable mTable = new MTable();
            // 1 row
            Row row = new Row();
            row.addItem("TE', 1");
            row.addItem("TE', 1");
            row.addItem("");
            row.addItem("");
            row.addItem("");
            row.addItem("");
            mTable.setRow(row);
            row = new Row();
            row.addItem("");
            row.addItem("");
            row.addItem("E, 3");
            row.addItem("+TE', 2");
            row.addItem("");
            row.addItem("E, 3");
            mTable.setRow(row);
            //      row.DebugRow();
            mTable.DebugTable();
        }
    }
} //конец для namespace

class MTable {
    ArrayList M = new ArrayList(); //
```

```

        public void setRow(Row row) { M.Add(row); }
        public void DebugTable() {
            //foreach (row r in M) {
                for(int i=0;i<M.Count;i++){
                    ((Row) M[i]).DebugRow();
                }
            //Console.WriteLine(r);
        }
    }

    class Row { // строка таблицы любого размера из
string - это колонка
        ArrayList row = new ArrayList(); //
        public void addItem (string item) {
            row.Add(item);
        }
        //        public ArrayList getRow() { return row; }
        public void DebugRow() {
            string[] str = null;
            foreach (string r in row) {
                Console.WriteLine(r);

                str = r.Split(',');
                foreach (string l in str) {
                    Console.WriteLine(l);
                }
            }
        }
    }
}

```

2. Определить размеры управляющей таблицы М. В соответствии с шагами алгоритм 3.10 построить управляющую таблицу М (см. раздел 3.3. “Практическая работа 7.”)

3. Выделить цепочку, принадлежащую языку порождаемому грамматикой G и написать правила вывода для цепочки.

4. Для LL(k)-грамматики G определить является ли она *сильно* LL(k)-грамматикой.

5. *Оформить работу согласно шагам.*

Лабораторная работа №14-15

Построить управляющую таблицу М для LR(k)-грамматики, написать правило вывода выделенной строки (см. раздел 3.4.). Описать работу алгоритма LR(k) анализатора.

1. Изучить раздел 3.4. Построение управляющей таблицы М для LR(k)-грамматики $G = (T, V, P, S)$.

2. В соответствии с шагами алгоритм 3.12 построить управляющую таблицу М (см. раздел 1.4. “Практическая работа 8.”)

3. Выделить цепочку, принадлежащую языку порождаемому грамматикой G и написать правила вывода для цепочки.

4. Для выделенной цепочки показать работу LR(k)-анализатора в соответствии с шагами алгоритм 3.11.

5. *Оформить работу согласно указанным пунктам и используемым шагам алгоритмов.*

Лабораторная работа №16

Применить алгоритм типа “перенос-свертка” для заданной грамматики $G=(T, V, P, S)$. Описать работу алгоритма.

1. Изучить алгоритм типа “перенос-свертка”

2. В соответствии с шагами алгоритм 3.13 по шагам рассмотреть работу алгоритма (см. раздел 3.4. “Практическая работа 8.”)

3. Выделить цепочку, принадлежащую языку порождаемому грамматикой G и написать правила вывода для цепочки.

4. Для выделенной цепочки показать работу LR(k)-анализатора в соответствии с шагами алгоритм 3.13 .

5. *Оформить работу согласно указанным пунктам и используемым шагам алгоритма.*

Приложение В. Задания к лабораторным работам

Для лабораторных работ 1-3 определены варианты автоматных языков:

1. $L=\{0\omega_1+(01)^* \mid \omega_1 \in \{0,1\}^*\}$
2. $L=\{01-(10)^*+\omega_101 \mid \omega_1 \in \{0,1\}^*\}$
3. $L=\{0(00)^*+01\omega_1 \mid \omega_1 \in (0,1,2)^*\}$
4. $L=\{\omega_1\omega_21 \mid \omega_1 \in \{1,0\}^+, \omega_2 \in \{1,0\}^+\}$
5. $L=\{1\omega_11-(00)^* \mid \omega_1 \in \{1,0\}^*\}$
6. $L=\{\omega_1-0\omega_2-0+\omega_3 \mid \omega_1 \in \{0,1\}^+, \omega_2 \in \{0,1\}^+, \omega_3 \in \{0,1\}^+\}$
7. $L=\{(0+1)(01)^*+\omega_1 \mid \omega_1 \in \{0,1\}^+\}$
8. $L=\{(0+1)^*\omega_1\omega_2 \mid \omega_1 \in \{0,1,2\}^+, \omega_2 \in \{0,1\}^+\}$
9. $L=\{(01)^*-1-(01)^*+\omega_1 \mid \omega_1 \in \{0,1\}^+\}$
10. $L=\{1(01)^*0-1+\omega_1 \mid \omega_1 \in \{0,1\}^+\}$
11. $L=\{(0+1)+\omega_1+(01)^*0 \mid \omega_1 \in \{0,1\}^+\}$
12. $L=\{0(000)^*(0+1)\omega_1 \mid \omega_1 \in \{0,1\}^*\}$
13. $L=\{1(01)^*(01)\omega_1 \mid \omega_1 \in \{0,1\}^*\}$
14. $L=\{00\omega_1+(1)^* \mid \omega_1 \in \{0,1\}^*\}$
15. $L=\{0(01)^*+1\omega_10 \mid \omega_1 \in \{1,0\}^+\}$
16. $L=\{1\omega_11\omega_21 \mid \omega_1 \in \{0,1\}^+, \omega_2 \in \{0,1\}^+\}$
17. $L=\{011\omega_11(0)^* \mid \omega_1 \in (0,1)^+\}$
18. $L=\{\omega_10-\omega_2 \mid \omega_1 \in \{0,1\}^+, \omega_2 \in \{0,1\}^+\}$
19. $L=\{\omega_1(1)^*0\omega_2 \mid \omega_1 \in \{1,2\}^*, \omega_2 \in \{0,1\}^*\}$

20. $L = \{10^+ \omega_1 (10)^* \omega_2 \mid \omega_1 \in \{1, 2\}^*, \omega_2 \in \{0, 1\}^*\}$
 21. $L = \{10 \omega_1 0^{-1} \omega_2 \mid \omega_1 \in \{0, 1\}^+, \omega_2 \in \{1, 2\}^*\}$
 22. $L = \{1^{-1} \omega_1 01 \omega_2 \mid \omega_1 \in \{1, 2\}^+, \omega_2 \in \{1, 2\}^*\}$

Лабораторные работы 3-6

А). Устранить из грамматики $G = (T, V, P, S)$ бесполезные символы, где

1. $P = \{S \rightarrow b, S \rightarrow F, S \rightarrow cFB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, F \rightarrow Ca, C \rightarrow d\}$
2. $P = \{S \rightarrow b, S \rightarrow cAB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, C \rightarrow Ca, F \rightarrow d\}$
3. $P = \{S \rightarrow b, S \rightarrow BA, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, F \rightarrow Ca, C \rightarrow d\}$
4. $P = \{S \rightarrow cB, B \rightarrow cB, B \rightarrow cA, A \rightarrow Ab, C \rightarrow Ca, F \rightarrow d\}$
5. $P = \{S \rightarrow c, F \rightarrow A, S \rightarrow cAB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, C \rightarrow Ca, C \rightarrow d\}$
6. $P = \{S \rightarrow cFCB, A \rightarrow ACb, A \rightarrow cC, B \rightarrow cB, C \rightarrow Ca, F \rightarrow d\}$
7. $P = \{S \rightarrow b, S \rightarrow CF, S \rightarrow cCB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, C \rightarrow Ca, F \rightarrow d\}$

В). Устранить из грамматики $G = (T, V, P, S)$ ϵ -правила, где

1. $P = \{S \rightarrow AB, A \rightarrow SA, A \rightarrow BB, A \rightarrow bB, A \rightarrow c, B \rightarrow c, B \rightarrow \epsilon\}$
2. $P = \{S \rightarrow b, S \rightarrow cAB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, B \rightarrow \epsilon\}$
3. $P = \{S \rightarrow bA, S \rightarrow bA, A \rightarrow Ab, A \rightarrow \epsilon, B \rightarrow cB, B \rightarrow \epsilon\}$
4. $P = \{S \rightarrow cB, B \rightarrow cB, B \rightarrow cA, A \rightarrow ACb, A \rightarrow \epsilon, C \rightarrow Ca, C \rightarrow \epsilon\}$
5. $P = \{S \rightarrow c, S \rightarrow cAB, S \rightarrow \epsilon, A \rightarrow c, B \rightarrow cB, B \rightarrow \epsilon\}$
6. $P = \{S \rightarrow cFCB, A \rightarrow ACb, A \rightarrow cC, B \rightarrow cB, B \rightarrow \epsilon, C \rightarrow Ca, F \rightarrow \epsilon\}$
7. $P = \{S \rightarrow b, S \rightarrow C, S \rightarrow cCB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, C \rightarrow Ca, C \rightarrow \epsilon\}$

С). Устранить из КС грамматики G цепные правила, где

1. $P = \{S \rightarrow AB, A \rightarrow S, A \rightarrow B, A \rightarrow bB, A \rightarrow c, B \rightarrow c\}$
2. $P = \{S \rightarrow b, S \rightarrow cAB, A \rightarrow Ab, A \rightarrow B, B \rightarrow cB, B \rightarrow b\}$
3. $P = \{S \rightarrow bA, S \rightarrow bA, A \rightarrow Ab, A \rightarrow S, B \rightarrow cB, B \rightarrow c\}$
4. $P = \{S \rightarrow cB, B \rightarrow cB, B \rightarrow cA, A \rightarrow C, A \rightarrow aB, C \rightarrow Ca, C \rightarrow cf\}$
5. $P = \{S \rightarrow c, S \rightarrow cAB, S \rightarrow a, A \rightarrow B, B \rightarrow cB, B \rightarrow f\}$
6. $P = \{S \rightarrow cFCB, A \rightarrow ACb, A \rightarrow C, B \rightarrow cB, B \rightarrow b, C \rightarrow Ca, F \rightarrow c\}$
7. $P = \{S \rightarrow b, S \rightarrow C, S \rightarrow cCB, A \rightarrow Ab, A \rightarrow C, B \rightarrow cB, C \rightarrow Ca, C \rightarrow b\}$

Д). Исключить левую рекурсию из КС – грамматики G , где

1. $P = \{S \rightarrow Bc, S \rightarrow Ad, A \rightarrow Sa, A \rightarrow AbB, A \rightarrow c, B \rightarrow Sc, B \rightarrow b\}$
2. $P = \{S \rightarrow FA, S \rightarrow c, A \rightarrow FS, A \rightarrow Sa, B \rightarrow SB, B \rightarrow b, F \rightarrow f\}$
3. $P = \{S \rightarrow Bb, B \rightarrow Sa, B \rightarrow cB, B \rightarrow Ac, A \rightarrow cSB, A \rightarrow a\}$
4. $P = \{S \rightarrow Ba, S \rightarrow Ab, A \rightarrow Sa, A \rightarrow AAb, A \rightarrow c, B \rightarrow Sb, B \rightarrow b\}$
5. $P = \{S \rightarrow ScB, S \rightarrow cAB, S \rightarrow c, A \rightarrow AbB, B \rightarrow b, B \rightarrow aA\}$
6. $P = \{S \rightarrow cFCB, A \rightarrow ACb, B \rightarrow cB, B \rightarrow b, C \rightarrow Ca, C \rightarrow c, F \rightarrow f\}$
7. $P = \{S \rightarrow AB, S \rightarrow SC, A \rightarrow BB, A \rightarrow Ab, A \rightarrow a, B \rightarrow b, C \rightarrow Ca, C \rightarrow b\}$

Лабораторные работы 3-8 задана КС-грамматика $G = (T, V, P, S)$, где

1. $T = \{i, =, *, (,)\}, V = \{S, F, L\}, P = \{S \rightarrow F = L, S \rightarrow L, F \rightarrow (* L), F \rightarrow i, L \rightarrow F\}$
2. $T = \{i, \&, *, (,)\}, V = \{S, F, L\}, P = \{S \rightarrow F \& L, S \rightarrow (S), F \rightarrow * L, F \rightarrow i, L \rightarrow F\}$
3. $T = \{i, ^\wedge, -, (,)\}, V = \{S, F, L\}, P = \{S \rightarrow (F ^\wedge L), F \rightarrow - L, F \rightarrow i, L \rightarrow F\}$

4. $T = \{i, *, -, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F)^* L, F \rightarrow - L, F \rightarrow i, L \rightarrow F\}$
5. $T = \{i, +, -, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F)+(L), F \rightarrow - L, F \rightarrow i, L \rightarrow F\}$
6. $T = \{i, +, -, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F)+(L), F \rightarrow - L, F \rightarrow i, L \rightarrow F\}$
7. $T = \{i, +, -, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F+L), S \rightarrow (F), F \rightarrow - L, F \rightarrow i, L \rightarrow F\}$
8. $T = \{i, \&, ^, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow F^{\wedge} L, S \rightarrow (F), F \rightarrow \& L, F \rightarrow i, L \rightarrow F\}$
9. $T = \{i, \&, ^, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow F^{\wedge} L, S \rightarrow (S), F \rightarrow \& L, F \rightarrow i, L \rightarrow F\}$
10. $T = \{i, \&, ^, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F^{\wedge} L), S \rightarrow (S), F \rightarrow \& L, F \rightarrow i, L \rightarrow F\}$
11. $T = \{i, \&, ^, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F^{\wedge} L), F \rightarrow \& L, F \rightarrow i, L \rightarrow F\}$
12. $T = \{i, \&, ^, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow \& F^{\wedge}, S \rightarrow (L), F \rightarrow L, F \rightarrow i, L \rightarrow F\}$
13. $T = \{i, *, :, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow F: L, S \rightarrow (L), F \rightarrow L^*, F \rightarrow i, L \rightarrow F\}$
14. $T = \{i, *, :, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F: L), F \rightarrow L^*, F \rightarrow i, L \rightarrow F\}$
15. $T = \{i, *, :, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F: L), S \rightarrow (F), F \rightarrow L^*, F \rightarrow i, L \rightarrow F\}$
16. $T = \{i, *, +, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F+ L), F \rightarrow L^*, F \rightarrow i, L \rightarrow F\}$
17. $T = \{i, *, +, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F+ L), F \rightarrow (L^*), F \rightarrow i, L \rightarrow F\}$
18. $T = \{i, *, +, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow F+ L, F \rightarrow (L^*), F \rightarrow i, L \rightarrow F\}$
19. $T = \{i, *, +, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow F+ L, S \rightarrow (S), F \rightarrow L^*, F \rightarrow i, L \rightarrow F\}$
20. $T = \{i, @, \&, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow (F@L), S \rightarrow (F\&L), F \rightarrow i, L \rightarrow F\}$
21. $T = \{i, +, -, (,)\}$, $V = \{S, F, L\}$, $P = \{S \rightarrow F+L, S \rightarrow (S), S \rightarrow (L-), F \rightarrow i, L \rightarrow F\}$

Лабораторные работы 14-16, правила G грамматики рассмотреть как правила LL(k) грамматики, а для работы 7 как правила LR(k).

Заключение

Представленный теоретический материал служит основой для практического проектирования грамматик по заданным языкам и эквивалентным грамматикам автоматов, для реализации синтаксических анализаторов на основе объектно-ориентированного подхода.

Лабораторные работы могут быть использованы для формирования курсовых и дипломных работ, а также исследовательских работ по теме теория автоматов и языков.

Библиографический список

1. Ахо А., Ульман Д. Теория синтаксического анализа перевода и компиляции // Пер. с англ. - М.: Мир, 1978.
2. Р.Хантер. Проектирование и конструирование компиляторов. Пер. с англ. – М.: Финансы и статистика, 1984.
3. Н.Вирт. Алгоритмы + структуры данных = программы. Пер. с англ. – М.: МИР, 1985.
4. C.Moore. Dynamical Recognizers: Real-time Language Recognition by Analog Computers. Theoretical Computer Science **201**, 1998, pp. 99-136.
5. W.Tabor. Fractal Encoding of Context Free Grammars in Connectionist Networks. University of Connecticut Expert Systems 17(1), 2000, pp. 41-56
6. А.С.Семенов. Построение класса фрактальных систем по шаблону на примере дерева Фибоначчи // Изв.РАН. Информационные технологии и Вычислительные системы. – М.: N2, 2005 стр.10-17.

7. *А.С.Семенов* Класс фрактальных автоматов для распознавания КС-языков // Программирование (рукопись).