

А.С. СЕМЕНОВ

ПОРОЖДАЮЩИЕ И РАСПОЗНАЮЩИЕ СИСТЕМЫ ФОРМАЛЬНЫХ ЯЗЫКОВ

С ПРИМЕРАМИ НА ЯЗЫКЕ ПРОГРАММИРОВАНИЯ С#

Москва 2020

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 форме ..	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	22
3.7. LR(k)-грамматики и пример построения LR(k)-анализатора	33
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. Аттрибутные транслирующие грамматики	
5.5. Методика разработки описания перевода	
5.6. Пример разработки АТ грамматики	
Глава 6. Разработка и реализация синтаксически управляемого перевода
6.1. L-аттрибутивные и S-аттрибутивные транслирующие грамматика	
6.2. Форма простого присваивания	
6.3. Аттрибутивный перевод для LL(1) грамматик	
6.4. S-аттрибутивный ДМП-процессор	
Приложение А. Порядок выполнения задач	
Приложение В. Задания для усвоения материала	
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 C# [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 C# [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 C#.

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

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.

ЧЕРНОВИК

Глава 1. Классификация грамматик и автоматов

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

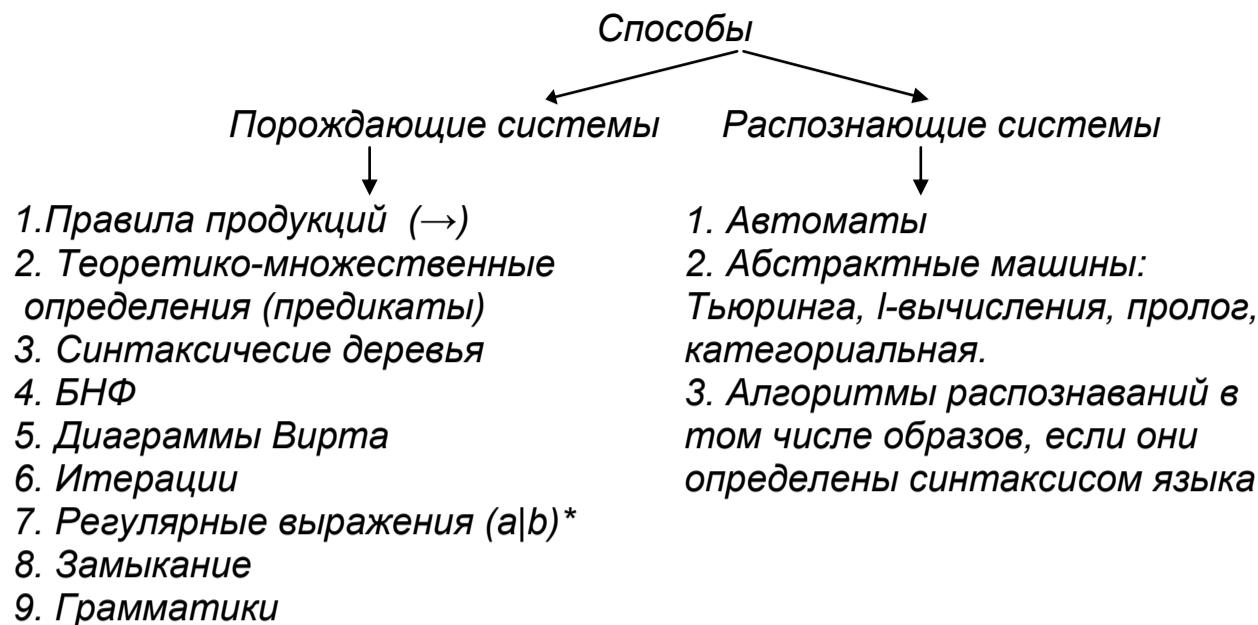


Рис. 1.1. Способы описания синтаксиса языков

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

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

то B будет множеством $\{2,4,8\}$ задается с помощью предиката, в данном случае это “ $x \subseteq N, x \text{ есть четное число}$ ”.

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 синтаксису, введем правило замены, которое обозначается символом “ \rightarrow ”. Тогда порождающая грамматика, для рассматриваемого предложения, примет вид:

предложение \rightarrow подлежащее, группа сказуемого

подлежащее \rightarrow он

группа сказуемого \rightarrow сказуемое, дополнение

сказуемое \rightarrow смотрит

дополнение \rightarrow кино

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

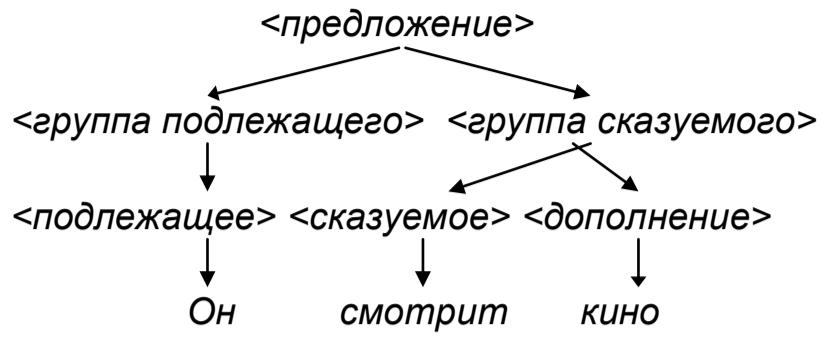


Рис. 1.2. Древовидная структура предложения

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 $\alpha ::= \beta | \beta\gamma$ можно представить $\alpha ::= \beta [\gamma]$.

БНФ для рассматриваемого примера.

$\text{<предложение>} ::= \text{<подлежащее>} \text{ <группа сказуемого>}$
 $\text{<подлежащее>} ::= \text{он} | \text{она}$
 $\text{<группа сказуемого>} ::= \text{<сказуемое>} \text{ <дополнение>}$
 $\text{<сказуемое>} ::= \text{смотрит} | \text{обожает}$
 $\text{<дополнение>} ::= \text{кино}$

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

$\text{<предложение>} \Rightarrow \text{<подлежащее>} \text{ <группа сказуемого>}$

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.

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

$\text{<подлежащее>} \text{ <группа сказуемого>} \Rightarrow \text{он} \text{ <группа сказуемого>}$

Полный вывод одного предложения будет таким:

$\text{<предложение>} \Rightarrow \text{<подлежащее>} \text{ <группа сказуемого>} \Rightarrow \text{он} \text{ <группа сказуемого>} \Rightarrow \text{он} \text{ <сказуемое>} \text{ <дополнение>} \Rightarrow \text{он} \text{ смотрит} \text{ <дополнение>} \Rightarrow \text{он} \text{ смотрит} \text{ кино}$

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

$\text{<предложение>} \Rightarrow^+ \text{ он} \text{ смотрит} \text{ кино}$

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: $T \rightarrow \{ F \} | (E)$

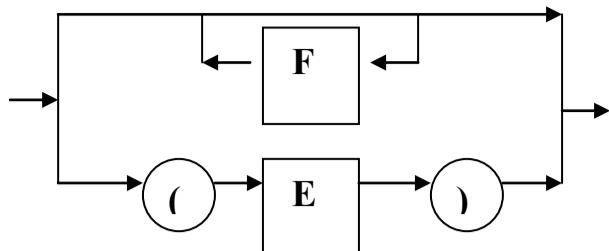


Рис. 1.3. Диаграммы Вирта

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 a в алфавите V называется любая конечная последовательность символов этого алфавита.

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

Длина цепочки a – число составляющих ее символов, обозначается | a |. Например, если цепочка символов a = defgabck, то длина a равна 8, | a | = 8. Длина ε равна 0, | ε | = 0.

Замыкание Клини (или звезда Клини) в математической логике и информатике – унарная операция над множеством строк либо символов. Замыкание Клини множества V обозначается V^* . Широко применяется в регулярных выражениях, на примере которых было введено Стивеном Клини 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 $V^* = \{\epsilon, 0, 1, 00, 11, 01, 10, 000, 001, 011, \dots\}$.

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

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

∅ – 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 находиться 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:

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

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

Цепочка $\beta \subseteq (T \cup V)^*$ **выводима** из цепочки $\alpha \subseteq (T \cup V)^+$ 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 $\alpha \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$, где $\beta_i, i = 1, 2, \dots, n$ и называют **альтернативой** правил вывода из цепочки а.

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

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

- $p_1: \text{Предложение} \rightarrow \text{подлежащее, группа сказуемого}$
- $p_2: \text{подлежащее} \rightarrow \text{он}$
- $p_3: \text{группа сказуемого} \rightarrow \text{сказуемое, дополнение}$
- $p_4: \text{сказуемое} \rightarrow \text{смотрит}$
- $p_5: \text{дополнение} \rightarrow \text{кино}$

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

$S_0 \rightarrow cBd$

$cB \rightarrow ccBd$

$B \rightarrow \epsilon$

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 $L(G) = \{\alpha \subseteq T^* \mid S_0 \Rightarrow^* \alpha\}$.

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

Цепочка $\alpha \subseteq (T \cup V)^*$, 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 :

$S_0 \rightarrow 0A1$

$0A \rightarrow 00A1$

$A \rightarrow \epsilon$

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

Грамматики 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 $L(G_1) = L(G_2) = \{0^n 1^n \mid n \geq 0\}$.

Грамматики G_2 и G_3 называются *почти эквивалентными* $L(G_2) = L(G_3)$, если $L(G_2) = L(G_3) \cup \{\epsilon\}$. 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

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

Тогда грамматики G_2 и G_3 почти эквивалентны, так как $L(G_2) = \{0^n 1^n \mid n \geq 0\}$, 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 Р	Пример правил грамматики $G = (T, V, P, S_0)$	Пример языка	Автомат или абстрактная машина
0	Общего вида, не накладывается никаких ограничений	$S_0 \rightarrow CD, C \rightarrow 0CA, C \rightarrow 1CB, AD \rightarrow 0D, BD \rightarrow 1D, A0 \rightarrow 0A, A1 \rightarrow 1A, B0 \rightarrow 0B, B1 \rightarrow 1B, C \rightarrow \epsilon, D \rightarrow \epsilon$	$\{\omega \omega \in \{0,1\}^*, \text{ язык состоит из цепочек четной длины, из } 0 \text{ и } 1\}$	Машина Тьюринга
1	Контекстно-зависимая (К3) $\alpha \rightarrow \beta, \alpha \leq \beta $	$S_0 \rightarrow 0B0$ $B \rightarrow 1 0BC$ $C0 \rightarrow 100$ $C1 \rightarrow 1C$	$\{0^n1^n0^n n \geq 1\}$	Машина Тьюринга с конечной лентой
2	Контекстно-свободная (КС) $A \rightarrow \alpha$	$S_0 \rightarrow AS_0B AB$ $A \rightarrow 0$ $B \rightarrow 1$ LL(k) грамматика LR(k) грамматика	$\{0^n1^n n \geq 1\}$	с магазинной памятью (МП-автомат: LL-разбор, РМП-автомат: LR-разбор) Распознаватели: LL(1) LR(1) подкласс LR(1) грамматики предшествования
3	Регулярная $A \rightarrow aB$ $B \rightarrow b$	$S_0 \rightarrow 0 0A$ $A \rightarrow 0B$ $B \rightarrow 1 1A$	$\{0(01)^n n \geq 0\}$	Конечный автомат (КА)

Грамматики классифицируются по виду правил вывода Р. В правилах вывода Р 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) А-грамматика – это грамматика $G = (T, V, P, S_0)$, у которой правила порождения P имеют вид по классификации Хомского $A \rightarrow bB$ (праволинейное правило) или $B \rightarrow a$ (заключительное правило), где $A, B \in V$, $a, b \in T$. 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 $A \rightarrow Bb$ либо $B \rightarrow a$, где $A, B \in V$, $a, b \in T$.

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 автоматами, 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 $S_0 \rightarrow \epsilon \in P$, 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: $S_0 01010 \mid A 1010 \mid B 010 \mid A 10 \mid B 0 \mid q_f, \epsilon$

2.2. Способы задания конечных автоматов

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

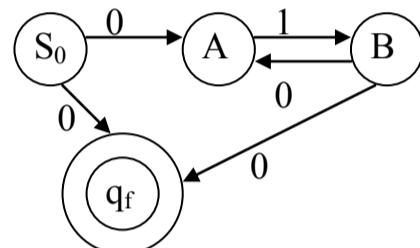


Рис. 2.2. Диаграмма переходов КА

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

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

Табл. 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 , отложенную 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 цепочки 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.

Формальное утверждение. Пусть L регулярный язык. Тогда существует целое $p \geq 1$ (длина накачки p принимается более длины самого длинной цепочки языка) 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.

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

Два способа:

1. построить распознаватель для заданного языка - конечный автомат;
2. выделить цепочку y^i для всех $i \geq 0$, $xy^i z \in L$, 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 рапознаватели для одних языков, построенных из других с помощью операций;
2. определить, что два различных автомата определяют один язык.

Теорема. Для каждого КС языка L , размер которого не является конечным, найдется такое натуральное число p , что для каждого $z \in L$ такое, что $|z| \geq p$, существуют такие $u, v, w, x, y \in \Sigma^*$, что $z = uvwxy$ с $|vwx| \leq p$ и $|vx| \geq 1$ такой, что для всех $i \geq 0$, $uv^ixy \in L$.

Эта лемма может быть переписана, как показано ниже, чтобы показать логический контекст. Обратите внимание, что в этой лемме p играет ту же роль, что и n в лемме накачки для обычных языков, а v и x , «место накачки», могут быть где угодно в z . Напомним, что «место накачки» v в лемме о накачке на обычном языке может находиться только в префиксе z длины n .

- (1) Для всех контекстно-свободных языков L ,
- (2) существует постоянная p такая, что
- (3) для всех $z \in L$ таких, что $|z| \geq p$,
- (4) существуют такие u, v, w, x, y и z , что
 - (i) $z = uvwxy$,
 - (ii) $|uwx| \leq p$,
 - (iii) $|ux| \geq 1$ и
 - (iv) для всех $i \geq 0$, $uv^ixy \in L$

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

Пример 2. Рассмотрим язык $L_{01} = \{0^n1^n \mid n \geq 1\}$, 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 состояния, 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.

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 xy — это первая часть w , состоящая из p последовательных экземпляров символа a .

Поскольку $|y| \geq 1$, она содержит по меньшей мере одну букву a , а $xy^2 z$ содержит больше букв a чем b . Следовательно, $xy^2 z$ 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.

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

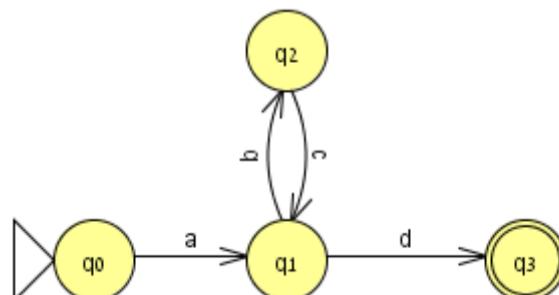


Рис. 2.3. Диаграмма переходов КА

Поскольку подстрока **bcbc** строки **abcabcd** 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 **abc**bcbc**cd** и **ad**.

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

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

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. Этую цепочку обозначим y по условиям леммы, и поскольку автомат примет строку как без части y , так и с повторяющейся частью y , условия леммы выполнены.

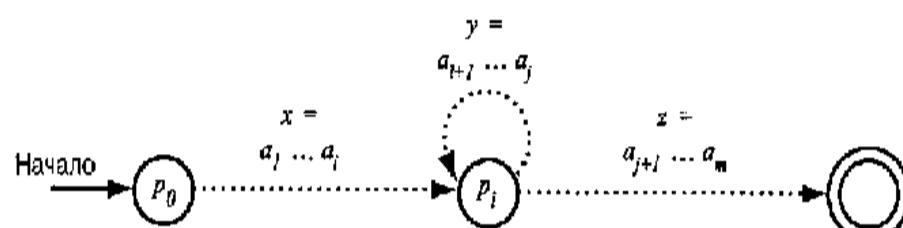
Доказательство. Пусть L – регулярный язык, тогда $L = L(\text{ДКА})$ для некоторого ДКА. Пусть ДКА имеет n состояний. Рассмотрим произвольную цепочку w длиной не менее n и скажем, $w = a_1a_2a_3..a_m$, где $m \geq n$ и каждый a_i есть входной символ. Для $i = 0, 1, 2, \dots, n$ определим состояние p_i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

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

Таким образом, x приводит в состояние p_i , y – 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 y не может быть пустой, поскольку i строго меньше j .



ис. 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 xy^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 y^k , k раз циклически проходит через p_j , и, наконец, по z переходит в допускающие состояния.

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

2.4. Способы преобразования НКА в ДКА

1. Построение НКА из нескольких автоматов.

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.

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

$\delta'(\{B\}, 0) = \{A, q_f\}$	$\delta'(\{S\}, 0) = \{A, q_f\}$	$\delta'(\{A\}, 1) = \{B\}$
$\delta'(\{S,B\}, 0) = \{A, q_f\}$	$\delta'(\{S,A\}, 0) = \{A, q_f\}$	$\delta'(\{S,A\}, 1) = \{B\}$
$\delta'(\{S,q_f\}, 0) = \{A, q_f\}$	$\delta'(\{B,A\}, 0) = \{A, q_f\}$	$\delta'(\{B,A\}, 1) = \{B\}$
$\delta'(\{B,q_f\}, 0) = \{A, q_f\}$	$\delta'(\{A,q_f\}, 1) = \{B\}$	$\delta'(\{S,B,A\}, 1) = \{B\}$
$\delta'(\{S,B,A\}, 0) = \{A, q_f\}$	$\delta'(\{S,B,q_f\}, 0) = \{A, q_f\}$	$\delta'(\{B,A,q_f\}, 0) = \{A, q_f\}$
$\delta'(\{B,A,q_f\}, 1) = \{B\}$	$\delta'(\{S,A,q_f\}, 0) = \{A, q_f\}$	$\delta'(\{S,A,q_f\}, 1) = \{B\}$
$\delta'(\{S,B,A,q_f\}, 0) = \{A, q_f\}$	$\delta'(\{S,B,A,q_f\}, 1) = \{B\}$	
- $F' = \{\{q_f\}, \{A, q_f\}, \{S, q_f\}, \{B, q_f\}, \{S, B, q_f\}, \{B, A, q_f\}, \{S, A, q_f\}, \{S, B, A, q_f\}\}$ Каждое из подмножеств F' является заключительным состоянием. При определении достижимого заключительного состояния выбирается соответствующее подмножество.
- $S = S_0$
- Достижимыми состояниями в ДКА являются $\{S\}$, $\{B\}$ и $\{A, q_f\}$, остальные состояния удаляются.

Табл. 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 = $(\{\{S\}, \{B\}, \{A, q_f\}\}, \{0, 1\}, \delta', S, \{\{A, q_f\}\})$, где

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

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

ДКА представляется следующим графом:

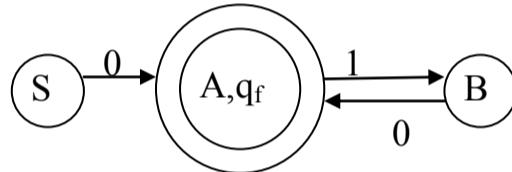


Рис. 2.5. Диаграмма переходов ДКА

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 $L = \{c\omega n \mid \omega \in \{+, -, k, f\}^*\}$, где ω обозначает множество всех цепочек, составленных из символов $\{+, -, k, f\}^*$. 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 леволинейной, бесконечной. Пример вывода цепочки $S_0 \Rightarrow Q-n \Rightarrow Qf-n \Rightarrow Qkf-n \Rightarrow Qdkf-n \Rightarrow Q-dkf-n \Rightarrow Q+d-dkf-n \Rightarrow c+d-dkf-n$.

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.:

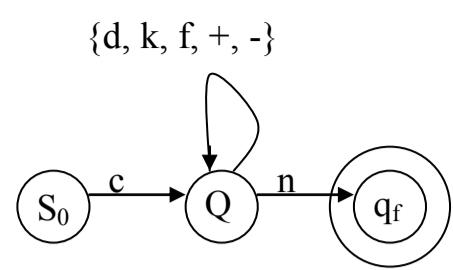


Рис. 2.6. Диаграмма переходов КА

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 $A \subseteq V$, $\alpha \subseteq (T \cup V)^*$.

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.

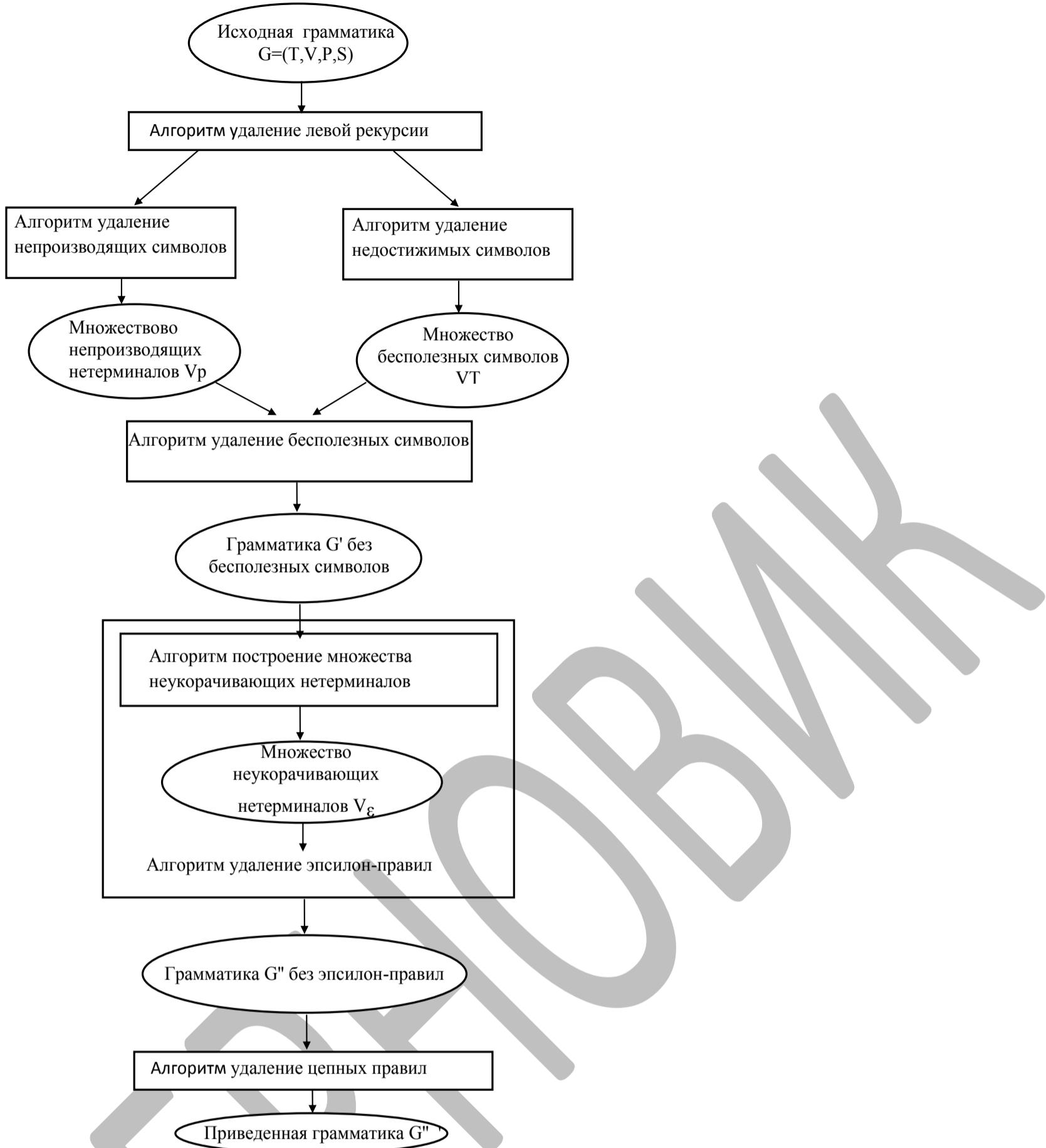


Рис.3.1 Метод преобразования грамматики к приведенной форме

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

При исходящем 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.

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 $A \rightarrow \epsilon$, и каждое правило из P отличное от $S \rightarrow \epsilon$, имеет вид $A \rightarrow a\alpha$, где $a \in T$, $\alpha \subset V^*$.

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. $S \rightarrow \epsilon$, если $\epsilon \in L(G)$, причем S не встречается в правых частях правил.

Алгоритм удаление левой рекурсии.

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

Выход: Эквивалентная приведенная КС грамматика G' .

Шаг 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 $i=1$.

```

 $V' := V$ 
 $P' := \emptyset$ 
foreach  $A \in V$ 
    if  $\exists (A \rightarrow A\alpha) \in P$ 
         $V_{rec} := V_{rec} \cup A$ 
         $V' := V' \cup A'$ 
    end if
end foreach

```

Шаг 2. Пусть множество A_i правил – это $A_i \rightarrow A_i\alpha_1 | \dots | A_i\alpha_m | \beta_1 | \dots | \beta_p$, где ни одна цепочка β_j не начинается с A_k , если $k \leq i$. Заменим A_i - правила правилами:

$$A_i \rightarrow \beta_1 | \dots | \beta_p | \beta_1 A'_i | \dots | \beta_p A'_i$$

$$A'_i \rightarrow \alpha_1 | \dots | \alpha_m | \alpha_1 A'_i | \dots | \alpha_m A'_i$$

где A'_i – новый символ. Правые части всех A_i - правил начинаются теперь с терминала или с A_k для некоторого $k > i$.

```
foreach ( $A \rightarrow \alpha_1 \dots \alpha_k \dots \alpha_n, \alpha_k \in T \cup V$ ) in  $P$ 
```

```

    if ( $A \in V_{rec}$ )
        if ( $\alpha_1 \neq A, n > 1$ )
             $P' := \{(A \rightarrow \alpha_1 \dots \alpha_k \dots \alpha_n A'), (A \rightarrow \alpha_1 \dots \alpha_k \dots \alpha_n)\} \cup P'$ 
        else
             $P' := \{(A \rightarrow \alpha_2 \dots \alpha_k \dots \alpha_n A'), (A \rightarrow \alpha_2 \dots \alpha_k \dots \alpha_n)\} \cup P'$ 
        end if
    else
        Заменить каждое правило вида  $A_i \rightarrow A_j\alpha$  правилами  $A_i \rightarrow \beta_1\alpha | \dots | \beta_m\alpha$ , где  $A_j \rightarrow \beta_1 | \dots | \beta_m$  – все  $A_j$ - правила.
         $P' := \{(A \rightarrow \alpha_1 \dots \alpha_k \dots \alpha_n) \cup P$ 
    end if
end foreach

```

Пример. Алгоритм удаления левой рекурсии из КС – грамматики $G = (\{v, +, *, (,)\}, \{S, F, L\}, P, S)$, где P состоит из правил

$$S \rightarrow S + F \mid F$$

$$F \rightarrow F * L \mid L$$

$$L \rightarrow v \mid (S)$$

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. $i = 1$ и $j = 2$.

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

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. $i = 2, j = 2$.

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

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

Обошли все нетерминалы

Получили правила новой грамматики G' :

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

Алгоритм удаление непроизводящих символов. Определение множества производящих нетерминальных символов V_p .

Если все символы цепочки из правой части правила вывода являются производящими, то нетерминал в левой части правила вывода также должен быть производящим: $V_p = \{A \mid A \Rightarrow^+ \alpha, A \in V, \alpha \in T^+\}$

Вход: КС $G = (T, V, P, S)$

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

```

 $V_p^0 := \emptyset$ 
 $V_p^1 := \emptyset$ 
do

```

```

 $V_p := V_p^i$ 
foreach ( $A \rightarrow \alpha_1 \dots \alpha_k \dots \alpha_n, \alpha_k \in T \cup V$ ) in P
    if ( $\alpha_k \in V_p \cup T$ ,  $\forall k \leq n$ ) // если есть производящее правило
         $V_p := V_p^i \cup \{A\}$  // положить  $V_p^i = \{A \mid (A \rightarrow \alpha) \in P, \alpha \in (T \cup V_{p^{i-1}})^+\}$ 
    end if
end foreach
while ( $V_p \neq V_{p^{i-1}}$ )

```

Алгоритм удаление недостижимых символов. Определение множества достижимых символов V_r . Если нетерминал в левой части правила грамматика является достижимым, то достижимы и все символы правой части этого правила.

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

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

```

 $V_r^0 := \{S\}$ 
 $V_r^{i-1} := \emptyset$  // дополнительное множество
 $i = 0;$ 
do
     $i = i + 1;$ 
     $V_r^i = V_r^{i-1}$ 
    foreach ( $A \in V$ )
        if ( $A \rightarrow \alpha X \beta \in P$ ) // Если нетерминал X достижим
             $V_r^i := V_r^i \cup X$  // Положить  $V_r^i = V_r^{i-1} \cup \{X \mid A \Rightarrow \alpha X \beta, X \in V, \alpha \in (V \cup T)^*, \beta \in (V \cup T)^* \text{ и } A \in V_r^{i-1}\}$ 
            if ( $\alpha \notin V_r^i$ )
                 $V_r^i := V_r^i \cup \alpha$ 
            end if
            if ( $\beta \notin V_r^i$ )
                 $V_r^i := V_r^i \cup \beta$ 
            end if
        end if
    end foreach
while ( $V_r^i \neq V_r^{i-1}$ )

```

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

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

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

Выполнить **Алгоритм удаления непроизводящих символов** и **Алгоритм удаление недостижимых символов** (см. выше).

$P' := \emptyset$

```

foreach ( $A \rightarrow \alpha X \beta \in P$ ) // для каждого правила  $\in P$ 
    if  $X \in V_p$  // X – производящий нетерминал
         $P' := P' \cup (A \rightarrow \alpha X \beta)$  // Добавляем правило, из которого выводится X
    end if
end foreach
 $T' := T \cap V_r$ 
 $V' := V \cap V_r$ 

```

Пример: УстраниТЬ из грамматики G бесполезные символы. $G = (T, V, P, S)$, где $V = \{A, B, C, S\}$, $T = \{a, b, c\}$, $P = \{S \rightarrow aC, S \rightarrow A, A \rightarrow cAB, B \rightarrow b, C \rightarrow a\}$.

множество $V_p = \{B, C, S\}$ А – непроизводящий символ

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

$V_r = \{C, S\}$ B – недостижим, так как A – непроизводящий символ

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

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

В грамматике с правилами вида $A \rightarrow \epsilon$ длина выводимой цепочки при переходе от k -го шага к $(k+1)$ -му уменьшается. Поэтому грамматики с правилами вида $A \rightarrow \epsilon$ называются *укорачивающими*.

Восходящий синтаксический разбор в укорачивающих грамматиках сложнее по сравнению с разбором в неукорачивающих грамматиках, т.к. при редукции необходимо отыскать такой фрагмент входной цепочки, в которую можно вставить пустой символ.

Для произвольной КС-грамматики, порождающей язык без пустой цепочки, можно построить эквивалентную неукорачивающую КС-грамматику.

Построение множества V_ϵ укорачивающих нетерминалов.

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

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

Выход: $V_\epsilon = \{A \mid A \xrightarrow{+} \epsilon, A \in V\}$.

```

 $V_\epsilon^0 = \emptyset$ 
foreach ( $A \in V$ )
    if ( $A \rightarrow \epsilon$ )  $\in P$  //  $A \rightarrow \epsilon$  правило укорачивающиеся
         $V_\epsilon := V_\epsilon \cup A$  // Положить  $V_\epsilon = V_\epsilon \cup \{A \mid (A \rightarrow \alpha) \in P \text{ и } \alpha \in (V_\epsilon)^+\}$ 
    end if
end foreach

```

```

 $V_\epsilon^1 := V_\epsilon$ 
 $V_\epsilon^0 := \emptyset$ 
 $i := 0$ 
do
     $i = i + 1$ 
     $V_\epsilon^{i-1} := V_\epsilon^i$ 
    foreach ( $A \rightarrow (C_1 \dots C_n), C_k \in V$ )  $\in P$  // для каждого правила  $\in P$ 
        if ( $C_k \in V_\epsilon, \forall k \leq n$ )
             $V_\epsilon^i := V_\epsilon^{i-1} \cup A$  // Положить  $V_\epsilon^i = V_\epsilon^{i-1} \cup \{A \mid (A \rightarrow \alpha) \in P \text{ и } \alpha \in (V_\epsilon^{i-1})^+\}$ .
        end if
    end foreach
    while ( $V_\epsilon^{i-1} \neq V_\epsilon^i$ )

```

Алгоритм удаление эпсилон-правил. Преобразование КС-грамматики с ϵ -правилами в эквивалентную НКС-грамматику.

Вход: КС $G = (T, V, P, S)$, V_ϵ

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

```

 $P' := \emptyset$ 
foreach ( $A \rightarrow \alpha_0 B_1 \alpha_1 B_2 \alpha_2 \dots B_k \alpha_k$ )  $\in P$  // для каждого правила  $\in P$ 
    // Если в правой части нетерминал, из которого выводится  $\epsilon$ -правило
    foreach ( $B_k \in V_\epsilon$ )
         $A := \text{delete}(B_k \text{ from } A)$ 
         $P' := P' \cup A$  // включить в  $P'$ , все правила вида  $A \rightarrow \alpha_0 X_1 \alpha_1 X_2 \alpha_2 \dots X_k \alpha_k$ , где  $X_j = B_j$  или  $X = \epsilon$ .
    end foreach
end foreach
if ( $S \in V_\epsilon$ )
     $V' = V \cup \{S'\}$ 
     $P' := P' \cup (S' \rightarrow S \mid \epsilon)$ 
end if

```

Пример: УстраниТЬ из грамматика G ϵ -правила. Преобразовать грамматику $G = (T, V, P, S)$, где $V = \{A, S\}$, $T = \{b, c\}$, $P = \{S \rightarrow cA, S \rightarrow \epsilon, A \rightarrow cA, A \rightarrow bA, A \rightarrow \epsilon\}$, в эквивалентную НКС-грамматику.

Шаг 1. Применяя алгоритм построения множества неукорачивающих нетерминалов, получаем $V_\epsilon^0 = \{S, A\}$, $V_\epsilon^1 = \{S, A\}$, значит $V_\epsilon^0 = V_\epsilon^1 = \{S, A\}$.

Шаг 2. Положить $P' = \emptyset$.

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

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

Алгоритм удаление цепных правил.

Правило вида $A \rightarrow B$, где A и B — нетерминалы называется цепным.

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

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

$P_c := \emptyset$ // множество цепных правил

$P_{nc} := \emptyset$ //множество нецепных правил

Шаг 1. Нахождение цепных и нецепных правил

```
foreach ( $A \rightarrow \alpha B \beta \in P$ )
    if ( $B \in V \text{ & } \alpha, \beta = \emptyset$ ) // Если  $A \rightarrow \alpha B \beta$  –цепное
         $P_c := P_c \cup (A \rightarrow B)$  // добавляем его в мн-во цепных  $P_c$ 
    else
         $P_{nc} := P_{nc} \cup (A \rightarrow \alpha B \beta)$  // иначе добавляем в мн-во нецепных  $P_{nc}$ 
    end if
end foreach
```

```
 $P' := \emptyset$ 
```

Шаг 2. Замена цепных на нецепные

```
foreach ( $R \rightarrow \alpha B \beta \in P$ )
    if  $B \rightarrow \gamma C \delta \in P_c$ 
         $P' := P' \cup (R \rightarrow \alpha \gamma C \delta \beta)$ 
    end if
end foreach
```

```
 $P' := P' \cup P_{nc}$ 
```

Пример УстраниТЬ из КС грамматики G цепные правила, $G = (\{S, F, L\}, \{v, +, *, (),\}, P, S)$, где P состоит из правил

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

Шаг 1. Получаем $V_S = \{S, F, L\}$, $V_F = \{F, L\}$, $V_L = \{L\}$.

Положить $P' = \emptyset$.

Шаг 2. Выберем первый нетерминал B из множества V. Множество правил, левая часть которых – нетерминальный символ R, правые части – это правые части нецепных правил исходной грамматики, в левой части которых находятся символы из множества V, 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 , в результате $G' = (T, V, P', S)$, где $P' = \{S \rightarrow S + F \mid F * L \mid (S) \mid v, F \rightarrow F * L \mid (S) \mid v, L \rightarrow v \mid (S)\}$

???Интегрированный алгоритм

Данные алгоритмы можно интегрировать, объединяя похожие части. Интегрируем алгоритмы 1-2:

```
 $VT_r := \{S_0\}$ 
 $V_p := \emptyset;$ 
 $P' := \emptyset$ 
foreach ( $A \in V$ )
    foreach ( $A \rightarrow \alpha X \beta \in P$ )
         $V_p := V_p \cup A$ 
         $VT_r := VT_r \cup X$ 
        if  $\alpha \notin VT_r$ 
             $VT_r := VT_r \cup \alpha$ 
        end if
        if  $\beta \notin VT_r$ 
             $VT_r := VT_r \cup \beta$ 
        end if
        if  $X \in V_p$ 
             $P' := P' \cup (A \rightarrow \alpha X \beta)$ 
        end if
    end foreach
end foreach
 $T' := T \cap VT_r$ 
 $V' := V \cap VT_r$ 
```

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.

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:

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

3.3. Свойства КС языков: лемма о накачке

Доказательство леммы накачки для КС языков.

Рассмотрим следующий КС языка в нормальной форме Хомского.

$S \rightarrow AB$ $A \rightarrow DE$ $B \rightarrow FD$ $E \rightarrow FG$
 $G \rightarrow DB \mid g$ $F \rightarrow HG$ $D \rightarrow d$ $H \rightarrow h$

Дерево вывода для строки $z = d(hd)^3hg(d^3dghdghd)$ приведено на рис.1. Обратите внимание на паттерн нетерминальной последовательности FGB-FGB-FGB, которая встречается на самом длинном стволе дерева синтаксического анализа. Он имеет повторяющийся шаблон FGB, который имеет листья "hd" слева и "d" справа. Легко видеть, что мы можем "вырастить" стебель, добавив столько повторяющихся шаблонов F-G-B, сколько мы хотим, или сократить его, удалив шаблон.

Полученное дерево по-прежнему является деревом синтаксического анализа грамматики. Из этого следует, что грамматика имеет дерево разбора для $z' = d(hd)^i hg(d^i dghdghd)$, для всех $i \geq 0$, что означает, что z' также находится на языке грамматики.

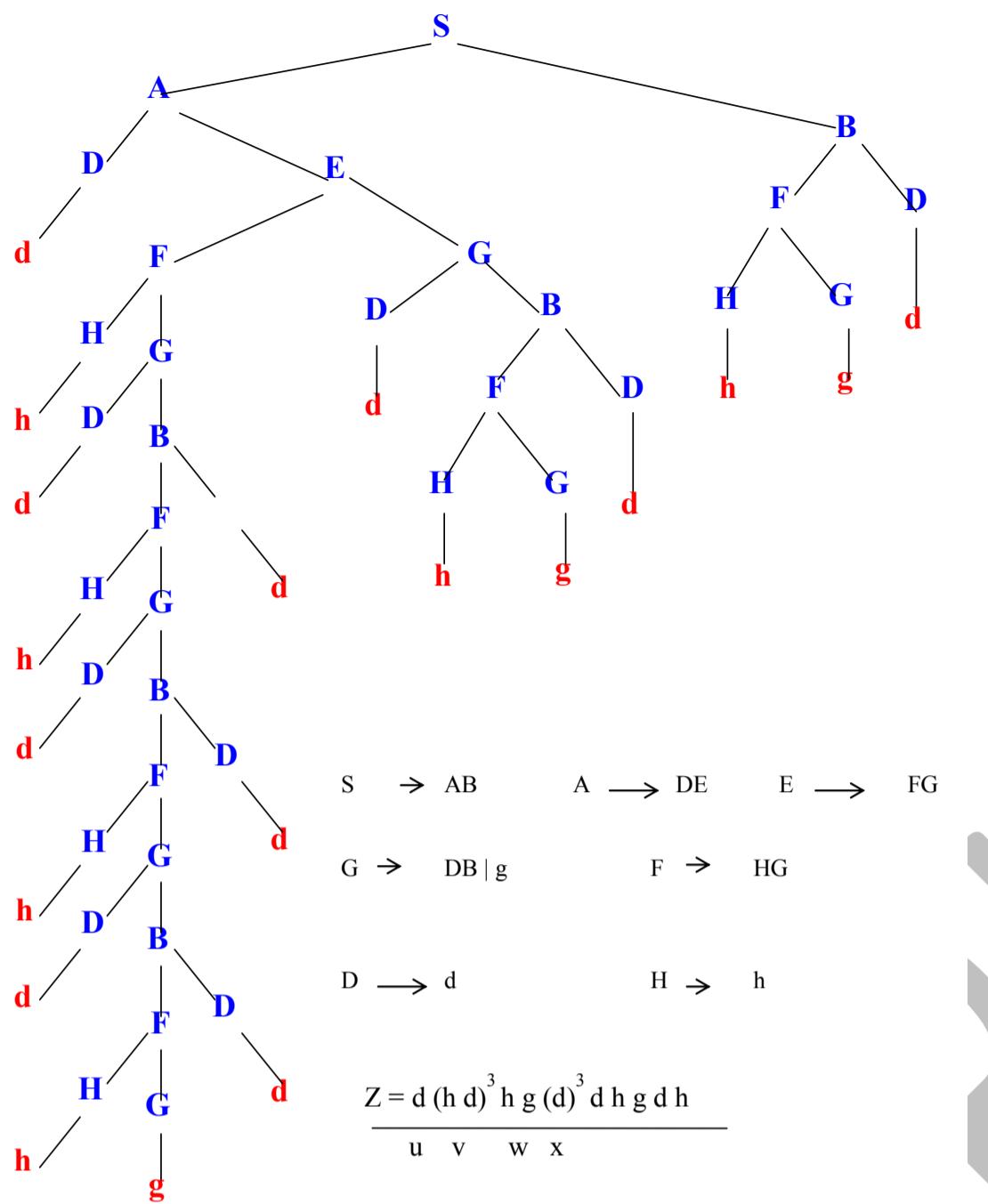


Рис. 2.6. Дерево разбора КС грамматики

uvwxyz
uvwxx
uvvwxx
uvwxyzxy

Реализация

Класс **UViWXiY** содержит поля **U,V,i,W,X,Y** соответствующие представлению строки $Z = UV^iWX^iY$

Класс **PumpingLema** следующие поля и методы:

- **private string Z** - цепочка символов, над которой будет производиться разбор
 - **private List<UViWXiY> ResultZ** - результат разбора Z, список структур, в которых содержится результат разбиения Z = **UViWXiY** при максимальном возможном значении i
 - **private List<string> WereChecked** - список подцепочек Z, для которых лема была проверена
 - **private void FindV()** - для цепочки Z перебираются все возможные подцепочки, которые принимаются за V
 - **private void FindVi(string V)** - для заданной подцепочки проверяется максимальное количество её непрерывных повторений. Если получается превысить полученное до этого число повторений i, результаты разбиения строки **UViWXiY**, результат заменяет предидущий.
 - **private void FindX()** - для цепочки Z перебираются все возможные подцепочки, которые принимаются за X
 - **private void FindXi(string X)** - для заданной подцепочки проверяется максимальное количество её непрерывных повторений. Если получается превысить полученное до этого число повторений i, результаты разбиения строки **UViWXiY**, результат заменяет предидущий.
 - **public void PrintResult()** - выводит на консоль результаты проверки

Цепочки, которые принимаются за V and X выбираются таким образом: во внешнем цикле изменяется длина, от 1 до половины длины Z (если длина V and X превысит половину длины Z , V and X не встретится более 1 раза), во внутреннем цикле выбирается r из Z , в котором будет начинаться получения подцепочки V and X длины $\text{Length } V$ and $\text{Length } X$ соответственно.

```

for (int LengthV = 1; LengthV <= Z.Length / 2; LengthV++)
{
    for (int p = 0; p < Z.Length - LengthV; ++p)
        {
            this.FindVi(Z.Substring(p, LengthV));
        }
}

for (int LengthX = 1; LengthX <= Z.Length / 2; LengthX++)
{
    for (int p = 0; p < Z.Length - LengthX; ++p)
        {
            this.FindXi(Z.Substring(p, LengthX));
        }
}

```

Для каждой выбранной подцепочки проверяется максимальное количество её непрерывных повторений.

В начале проверки проверяется, была ли просмотрена подпоследовательность, выбранная в качестве V. Если не была, то заносится в список **WereChecked**, алгоритм продолжается. Цепочку Z можно покрыть V.Length способами. Например Для V Длинной 3 Z можно покрыть 3 способами:

Для Каждого из этих способов не сложно подсчитать самое большое количество непрерывных повторений, при этом необходимо запоминать позицию начала соответствующей последовательности повторений. После прохода получившееся максимальное количество повторений сравнивается с предыдущим полученным, и если оно превосходит, предидущие результаты стираются, в стек помещаются новые. Если одному покрытию соответствует несколько одинаковых по количеству повторений последовательностей, в стек добавляются все эти варианты. После завершения проверки для каждого покрытия, в стеке будет содержаться информация разбиении $Z = UV^iWX^iY$, для заданной подцепочки V и максимально возможной i, если таких разбиений несколько - тоже. Результат из стека переносится в **ResultZ** только если значение i больше 1 и не меньше найденного ранее при другом выборе V.

Аналогично приводится пример с подцепочкой X

Пример. Программа корректно выявила разбиения строки $Z = UV^iWX^iY$ с наибольшим возможным значением i.

abcdefedcba

aaaaaaaa

a^7

ddddddd

d^7

aaabaaaadaaaa

aaabaaaad a^3

aaab a^3 daaa

a^3 baaadaaa

d^3

abbbbc

a b^4 c

abbbccddde

a b^3 cccddde

abbb c^3 ddde

abbbccc d^3 e

abcdabcdabcdabcd

$abcd^4$

aaabcdabcdabcdabcdaab

aa abcd 4 daab

aaabcdbcdbcdabcdaab

aaabcdcbcdaaabcbcdabcd a³ bcdbcdabcd
 aaabcdcbcdaabcbcdabcd a³ bcdbcdcbcdaaabcbcdabcd
 a³ bcdbcdcbcdaaabcbcd
 d³ dbcdaaabcbcdabcd
 aaabcdbc d³
 dbcdcaaabcdbcdbcdaaa bcd³
 aaabcdcbcdaaaa bcd³ aaabcdcbcdaabcbcdabcd
 aaa bcd³ aaabcdbc d³ dcdaaabcbcdabcd
 aaabcdcbcdaabcbcd³

abbbbabbabbabbabbabb
 abbbbabb d³ bbabbba b⁴
 abbbbabbba b⁴ abbbb
 abbbba b⁴ abbbbabbbb
 a b⁴ abbbbab d³ bbbabbbb
 abbbb⁴

abcdbcdbcdcbcd
 a bcd⁵
 abcbcbcd
 a bc³ bd
 ab cb³ d d³ cd

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

$$M\bar{P} = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$$

Q – конечное множество состояний устройства управления;

Σ - конечный алфавит входных символов;

Γ - конечный алфавит магазинных символов;

δ - функция переходов, отображает множества $Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma$ 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 $\omega = \epsilon$, то считается, что вся входная цепочка прочитана);

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

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

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

Если $a \neq \epsilon$, 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 $a = \epsilon$, то текущий входной символ в этом такте называется ϵ -таком, 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 таkt 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 $M = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$, если $(q_0, \omega, z_0) \vdash^* (q, \epsilon, \gamma)$ для некоторых $q \in F$ и $\gamma \in \Gamma^*$.

Язык распознаваемый МП-автоматом, называется множество цепочек:

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

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

Цепочка символов языка L : aaabbb

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

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

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

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

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

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

Последовательность тактов:

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

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

1. Если $A \rightarrow \alpha$ - правило грамматики G , то $\delta(q, \epsilon, A) = (q, \alpha)$.
2. $\delta(q, a, a) = \{(q, \epsilon)\}$ 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: $M = (\{q\}, \{v, +, *, (),\}, \{v, +, *, (),\}, \{S, F, L\}, \delta, q_0, S, \{q\})$, в котором функция переходов δ определяется следующим образом:

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

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

$$(q_0, v^*(v+v), S) \vdash_1 (q, v^*(v+v), F) \vdash_1 (q, v^*(v+v), F*L) \vdash_1 (q, v^*(v+v), L*L) \vdash_1 (q, v^*(v+v), v*L) \vdash_1 (q, v^*(v+v), v*v) \vdash_1 (q, (v+v), L) \vdash_1 (q, (v+v), (S)) \vdash_1 (q, (v+v), S) \vdash_1 (q, (v+v), S+F) \vdash_1 (q, (v+v), F+F) \vdash_1 (q, (v+v), L+F) \vdash_1 (q, (v+v), v+F) \vdash_1 (q, (+v), +F) \vdash_1 (q, v, F) \vdash_1 (q, v, L) \vdash_1 (q, v, v) \vdash_1 (q, v,)) \vdash_1 (q, \epsilon, \epsilon).$$

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 *находящим* (предсказывающим) анализатором.

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

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

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 $\text{РМП} = (\{q, r\}, \{v, +, *, (,), S, F, L, \perp\}, \delta, q, \perp, \{r\})$, где функция переходов δ определена следующим образом:

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

Б). При анализе входной цепочки $v+v^*v$ расширенный РМП-автомат выполнит следующую последовательность тиков:

- $(q, v+v^*v, \perp) \xrightarrow{1} (q, +v^*v, \perp v)$
 $\xrightarrow{2} (q, +v^*v, \perp L)$
 $\xrightarrow{3} (q, +v^*v, \perp F)$
 $\xrightarrow{4} (q, +v^*v, \perp S)$
 $\xrightarrow{5} (q, v^*v, \perp S+)$
 $\xrightarrow{6} (q, *v, \perp S+v)$
 $\xrightarrow{7} (q, *v, \perp S+L)$
 $\xrightarrow{8} (q, *v, \perp S+F)$
 $\xrightarrow{9} (q, v, \perp S+F^*)$
 $\xrightarrow{10} (q, \epsilon, \perp S+F^*v)$
 $\xrightarrow{11} (q, \epsilon, \perp S+F*L)$
 $\xrightarrow{12} (q, \epsilon, \perp S+F)$
 $\xrightarrow{13} (q, \epsilon, \perp S)$
 $\xrightarrow{14} (r, \epsilon, \epsilon)$

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, ..., 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 $\text{ДМП} = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$, 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 $\delta(q, \epsilon, z) = \emptyset$;
2. $\delta(q, a, z) = \emptyset$ для всех $a \in \Sigma$, и $\delta(q, \epsilon, z)$ 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 косвенной 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 $G = (\{a, b, c, -, +, ()\}, \{E, T, F, op\}, P, S)$, где $P = \{E \rightarrow T \text{ op } T, \text{ op} \rightarrow + | -, T \rightarrow F | (E), F \rightarrow a | b | c\}$

2. Шаги реализации.

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

Привести вывод выделенной цепочки:

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

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

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

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 A можно пройти, не читая вообще никаких входных символов, то такая “нулевая ветвь” должна помечаться всеми символами, которые могут следовать за A (это необходимо для принятия решения о переходе на эту ветвь).

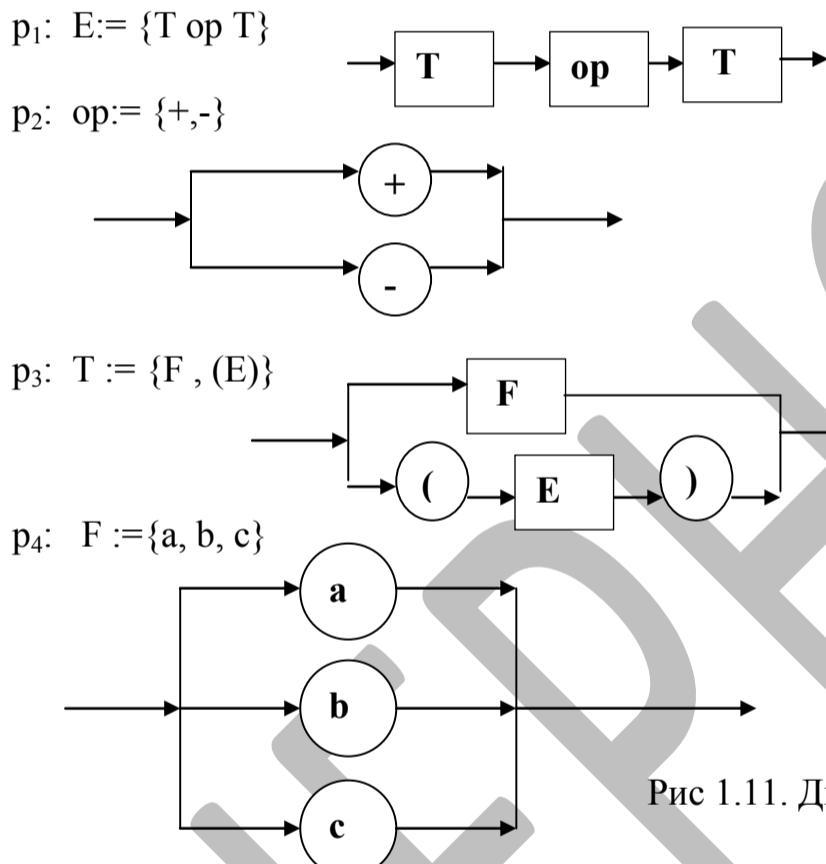
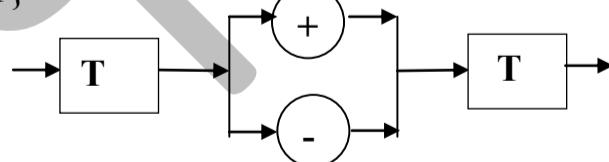


Рис 1.11. Диаграммы Вирта.

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 :

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



p₄: T := {{a, b, c}, (E)}

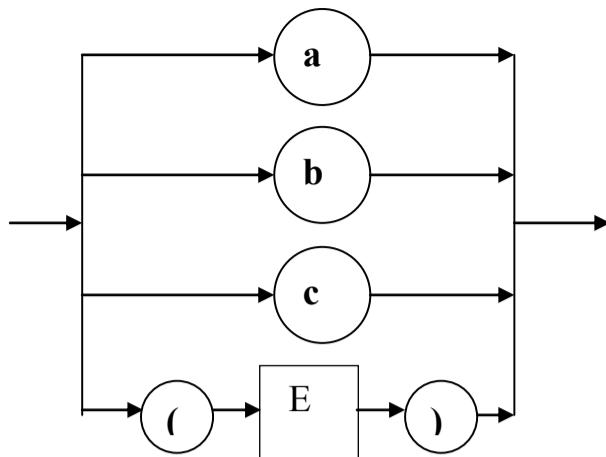


Рис 1.12. Преобразованные диаграммы Вирта.

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, порождаемому 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:

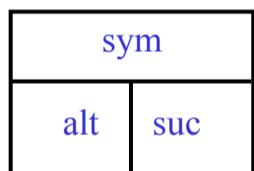


Рис. 1. 14. Структура узла Node

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.

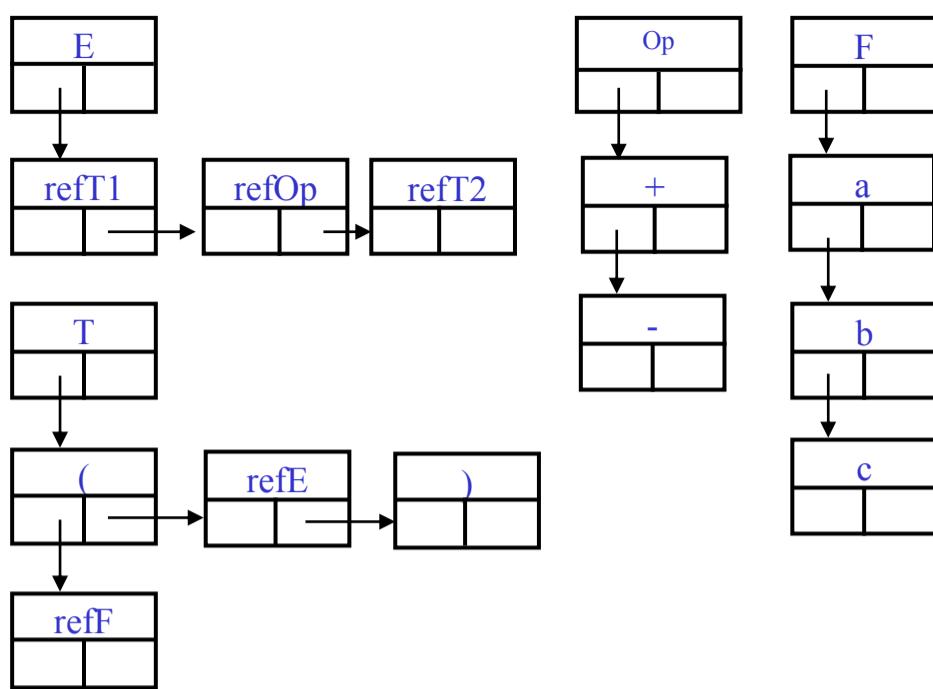


Рис. 1.15. Структура данных для диаграмм п.5

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

```

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
// контексто-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}"".", c, 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+f fg"; // error
    // string str = "a+(b-c)";
    Console.Write("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.

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

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

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

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

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.

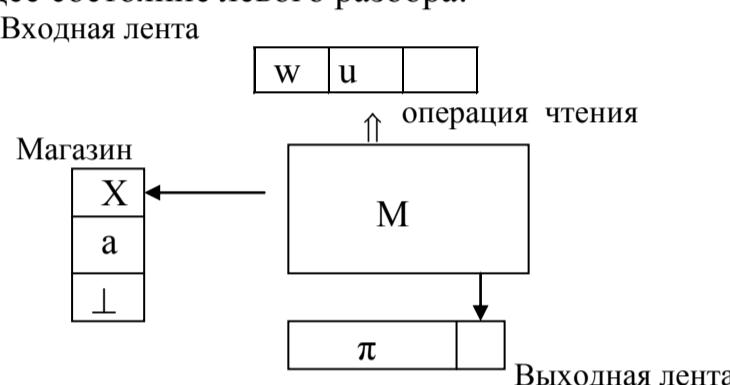


Рис. 1.18. LL(k)-анализатор

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 $(V_p \cup T \cup \{\perp\}) \times (T \cup \{\epsilon\})$ 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 $M(X, a) = (\beta, i)$ следующая:

1. $(x, Xa, \pi) \vdash (x, \beta a, \pi)$, если $M(X, a) = (\beta, i)$. 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. $(ax, a, \pi) \vdash (x, a, \pi)$, 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 (ϵ, \perp, π) , что соответствует элементу управляющей таблицы $M(\perp, \epsilon) = \text{ДОПУСК}$, то его работа прекращается, и выходная цепочка π является левым разбором входной цепочки.

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

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

Если $(\omega, S\perp, \epsilon) \vdash (\epsilon, \perp, \pi)$, 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 $(V \cup T \cup \{\perp\}) \times (T \cup \{\epsilon\})$ по правилам:

1. Если $A \rightarrow \beta$ – правило грамматики с номером i , то $M(A, a) = (\beta, i)$ для всех $a \neq \epsilon$, принадлежащих множеству $\text{FIRST}(A)$. Если $\epsilon \in \text{FIRST}(\beta)$, 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. $M(\perp, \epsilon) = \text{ДОПУСК}$.
4. В остальных случаях $M(X, a) = \text{ОШИБКА}$ для $X(V \cup T \cup \{\perp\})$ и $a \in T \setminus \{\epsilon\}$

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:

$$\text{FIRST}_k(\beta) \cap \text{FOLLOW}_k(A) \cap \text{FIRST}_k(\gamma) \cap \text{FOLLOW}_k(A) = \emptyset, \text{ то такая}$$

КС-грамматика называется *сильно LL(k)-грамматикой*.

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

Проектирование LL(k)-анализатора.

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

Б) 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 $(T \cup \{\epsilon\})$.

	i	$($	$)$	$+$	$*$	ϵ
S	$TE', 1$	$TE', 1$				
E'			$\epsilon, 3$	$+TE', 2$		$\epsilon, 3$
T	$PT', 4$	$PT', 4$				
T'			$\epsilon, 6$	$\epsilon, 6$	$*PT', 5$	$\epsilon, 6$
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 $\text{FIRST}(TE') = \{ (, i \}$, два терминальных символа, то:

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

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

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

$p_3: E' \rightarrow \epsilon$ 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 $\text{FOLLOW}(E') = \{) \}, \epsilon \}$.

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

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

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

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 $M(\perp, \epsilon)$ 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 $(i+i^*i, S\perp, \epsilon)$. Значение управляющей таблицы $M(S, i) = TE', 1$, 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 M
$(i+i^*i, TE' \perp, 1) \vdash$	$M(T, i) = PT', 4$
$(i+i^*i, PT'E' \perp, 14) \vdash$	$M(P, i) = i, 8$
$(i+i^*i, iT'E' \perp, 148) \vdash$	$M(i, i) = ВЫБРОС$

Шаг 3. Алгоритм находится в конфигурации $(i+i^*i, iT'E' \perp, 148)$. $M(i, i) = ВЫБРОС$. 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 M
$(+i^*i, T'E' \perp, 148) \vdash$	$M(T', +) = \epsilon, 6$
$(+i^*i, E' \perp, 1486) \vdash$	$M(E', +) = + TE', 2$
$(+i^*i, +TE' \perp, 14862) \vdash$	$M(+, +) = ВЫБРОС$
$(i^*i, TE' \perp, 14862) \vdash$	$M(T, i) = PT', 4$
$(i^*i, PT'E' \perp, 148624) \vdash$	$M(P, i) = i, 8$
$(i^*i, iT'E' \perp, 1486248) \vdash$	$M(i, i) = ВЫБРОС$
$(*i, T'E' \perp, 1486248) \vdash$	$M(T', *) = *PT', 5$
$(*i, *PT'E' \perp, 14862485) \vdash$	$M(*, *) = ВЫБРОС$
$(i, PT'E' \perp, 14862485) \vdash$	$M(P, i) = (i, 8)$
$(i, iT'E' \perp, 148624858) \vdash$	$M(i, i) = ВЫБРОС$
$(\epsilon, T'E' \perp, 148624858) \vdash$	$M(T', \epsilon) = (\epsilon, 6)$
$(\epsilon, E' \perp, 1486248586) \vdash$	$M(E', \epsilon) = (\epsilon, 3)$
$(\epsilon, \perp, 14862485863) \vdash$	

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

Так как значение $M(\perp, \epsilon) = ДОПУСК$, 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, называются **LR(k)-грамматиками** (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).

Для определения **LR(k)-грамматики** используются:

1. Множество $\text{FIRST}_k(\gamma)$, 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 $\text{FIRST}_k(\gamma)$. Формально:

$$\text{FIRST}_k(\gamma) = \{x \mid \gamma \Rightarrow_1^* xw \text{ и } |x| = k \text{ или } \gamma \Rightarrow_1^* x \text{ и } |x| < k\}$$

2. Пополненной грамматикой, полученной из КС-грамматики $G = (T, V, P, S)$, называется грамматики $G' = (V \cup \{S'\}, T, P = \{S' \rightarrow S\}, S')$. 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 грамматики 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:

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

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

для которых $\text{FIRST}_k(w) = \text{FIRST}_k(y)$ следует, что $\alpha Ay = \gamma Bx$.

То есть, если $\alpha \beta w$ и $\alpha \beta y$ – правовыводимые цепочки пополненной грамматики G' , у которых $\text{FIRST}_k(w) = \text{FIRST}_k(y)$ 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 $\alpha \beta y$ к αAy .

Поскольку правило грамматики $A \rightarrow \beta$ не зависит от w, то из определения LR(k)-грамматики следует, что во множестве $\text{FIRST}_k(w)$ содержится информация достаточная для определения основы. **Основа – это** 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.

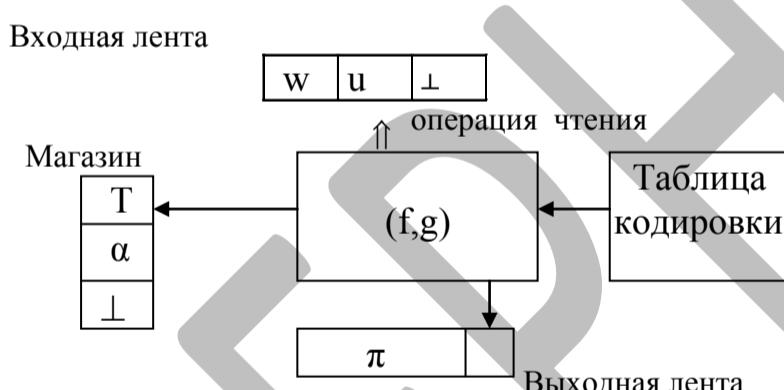


Рис. 1.18. LR(k)-анализатор

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

Два способа построения LR(k) анализаторов

Существуют два способа построения LR(k) анализаторов:

- 1) Определение активных префиксов - использование расширенного магазинного алфавита (алгоритм перенос-свертка)
- 2) На основе грамматических вхождений.

Построим двумя способами LR(k) анализатор для заданной грамматики $G = (T = \{i, \&, *, (,)\}, V = \{S, F, L\}, P = \{p_1: S \rightarrow F \& L, p_2: S \rightarrow (S), p_3: F \rightarrow * L, p_4: F \rightarrow i, p_5: L \rightarrow F\}, S\}$

1. Построение LR(k) анализатора способом использования расширенного магазинного алфавита

Способ использования расширенного магазина состоит из трех шагов:

Шаг 1. Определение активных префиксов;

Шаг 2. Построение управляющей таблицы;

Шаг 3. Применение алгоритма «перенос-свёртка».

Шаг 1. Определение активных префиксов

Для каждого магазинного символа (за исключением S_0 и \perp), кодируемая цепочка является префиксом правой части некоторого правила грамматики. Например, магазинный символ « $&_1$ », находящийся на вершине магазина, означает что:

1. В верхушке магазина находится терминальный символ грамматики « $&$ »

2. В верхней части магазина находится кодируемая этим символом цепочка символов « $F&$ », то есть префикс основы – правой части правила 1.

Индекс каждого символа соответствует номеру правила, префикс правой части которого кодируется этим символом. И наоборот, каждый непустой префикс правой части правила кодируется некоторым магазинным символом. Например, правая часть правила 3 имеет два непустых префикса: « $*$ » и « F », которые соответственно кодируются символами « $*_3$ » и « L_3 ».

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.

Таблица 1. Закодированные символы для заданной грамматики.

Символ грамматики	Магазинный Символ	Кодируемая цепочка
S	S_0	$\perp S$
	S_2	(S
F	F_1	F
	F_5	F
L	L_1	F&L
	L_3	*L
I	i_4	I
($(_2$	(
*	$*_3$	*
&	$&_1$	F&
)	$)_2$	(S)

Шаг 2. Построение управляющей таблицы;

Управление алгоритмом осуществляется при помощи двух функций, приведенных в таблице 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.

Работа алгоритма описывается в терминах конфигураций, представляющих собой тройки вида $(\alpha T, ux, \pi)$, где αT – цепочка магазинных символов (T -верхний символ магазина), ux – необработанная часть входной цепочки, π – выход (строка из номеров правил), построенный к настоящему моменту времени.

Таблица состоит из двух подтаблиц – функции действия и функции переходов. Входным символам с ленты соответствуют столбцы таблицы, символам магазина – строки.

Алгоритм 3.11 Построение управляющей таблицы для LR(1) грамматики

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

Выход: Управляющая таблица M для грамматики G .

Функция действий $f(u)$ определяется на множестве $(V_p \cup \{\perp\}) \times (T \cup \{\epsilon\})$ по правилам:

1. Если $A \rightarrow \beta$ – правило грамматики с номером i , то для конфигурации $(\alpha T, ux, \pi)$, где T кодирует цепочку β , $f(u) = C(i)$.
2. Если $A \rightarrow \beta$ – правило грамматики с номером i , то для конфигурации $(\alpha T, ux, \pi)$, где T кодирует некий префикс цепочки β (но не саму основу), $f(u) = \Pi$.
3. Для конфигурации (S_0, \perp, π) , где S_0 кодирует цепочку $\perp S$, $f(u) = \text{ДОПУСК}$.
4. В противном случае, $f(u) = \text{ОШИБКА}$.

Функция переходов $g(X)$ определяется на множестве $(V_p \cup \{\perp\}) \times (V \cup T \cup \{\epsilon\})$ по правилам:

- Если для конфигурации $(\alpha T, ax, \pi)$ для входного символа $a \in (V \cup T)$ в таблице 1 существует символ a_i , совместимый с цепочкой αTa , то $g(X) = a_i$.
- В противном случае, $g(X) = \text{ОШИБКА}$.

		функция действий $f(u)$					функция переходов $g(X)$							
	i	*	&	()	\perp	S	F	L	i	*	&	()
S_0	Π	Π	Π	Π	Π	Δ		F_1		i_4	$*_3$		$(_2)$	
S_2	Π	Π	Π	Π	Π			F_5		i_4	$*_3$		$(_2)_2$	
F_1	Π	Π	Π	Π	Π						$\&_1$			
F_5	$C(5)$	$C(5)$	$C(5)$	$C(5)$	$C(5)$									
L_1	$C(1)$	$C(1)$	$C(1)$	$C(1)$	$C(1)$									
L_3	$C(3)$	$C(3)$	$C(3)$	$C(3)$	$C(3)$									
i_4	$C(4)$	$C(4)$	$C(4)$	$C(4)$	$C(4)$									
$(_2)$	Π	Π	Π	Π	Π		S_2	F_1		i_4	$*_3$		$(_2)$	
$*_3$	Π	Π	Π	Π	Π			F_5	L_3	i_4	$*_3$			
$\&_1$	Π	Π	Π	Π	Π			F_5	L_1	i_4	$*_3$			
$)_2$	$C(2)$	$C(2)$	$C(2)$	$C(2)$	$C(2)$	$C(2)$								
\perp	Π	Π	Π	Π	Π		S_0	F_1		i_4	$*_3$		$(_2)$	

Шаг 3. Применение алгоритма «перенос-свёртка».

Алгоритм 3.12. LR(k)-алгоритм разбора

Вход: Анализируемая цепочка $z = t_1 t_2 \dots t_j \dots t_n \in T^*$, 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.

- $j := 0$.
- $j := j+1$. Если $j > n$, то выдать сообщение об ошибке и перейти к шагу 5.
- Определить цепочку и 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.
- Применить функцию действия f из строки таблицы M , отмеченной верхним символом магазина T , к цепочке u :
 - $f(u) = \text{ПЕРЕНОС } (\Pi)$. Определить функцию перехода $g(t_j)$ из строки таблицы M , отмеченной символом T из верхушки магазина. Если $g(t_j) = T'$ и $T' \in V_p \cup \{\perp\}$, 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 $T' \in V_p \cup \{\perp\}$, 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.
 - $f(u) = \text{ОШИБКА } (O)$. Выдать сообщение об ошибке и перейти к шагу 5.
 - $f(u) = \text{ДОПУСК } (\Delta)$. Объявить цепочку, записанную на входной ленте, правым разбором входной цепочки z .
- Останов.

Пример

Рассмотрим, например, строку таблицы, отмеченную закодированным символом $\&_1$, который представляет собой префикс правой части правила $p_1: S \rightarrow F \& L$. Пусть i – текущий входной символ. Согласно правилу p_1 за $\&_1$ может следовать нетерминал L_1 или символ, представляющий подцепочку, порождающую L_1 . Цепочка, порождаемая L_1 , может начинаться либо с i_4 (правило p_4), либо с F_5 (правило p_5), которое в свою очередь порождает цепочку, начинающуюся с $*_3$ или снова i_4 . Следовательно, для входного символа i в магазин можно втолкнуть только символ i_4 , который является кодированным представлением префикса правой части правила p_4 .

Строка таблицы, отмеченная маркером дна магазина (\perp), соответствует начальной конфигурации алгоритма. В первый момент времени в магазин можно записать символ S_0 или символы $F_1, i_4, *_3, (_1$, которые представляют собой префиксы цепочек, выводимых из начального вхождения S_0 (грамматического вхождения начального символа грамматики S , входящего в правую часть нулевого правила вывода пополненной грамматики).

Рассмотрим последовательность тиков LR(k)-алгоритма при анализе входной цепочки $(*i\&i)$:

Шаг 1. Начальная конфигурация алгоритма – $(\perp, (*i\&i)\perp, \varepsilon)$ (в магазине находится маркер дна магазина, текущий входной символ – «»).

Для строки управляющей таблицы, отмеченной символом \perp , $f(\) = \text{ПЕРЕНОС}$ и $g(\) = (_2$, поэтому:

- В магазин записывается символ $(_2)$.
- Входная головка сдвигается на один символ вправо.

Алгоритм переходит в конфигурацию $(\perp, (*i\&i)\perp, \varepsilon)$.

Шаг 2. Для строки таблицы, отмеченной символом $(_2, f(*) = \text{ПЕРЕНОС}$ и $g(*) = *_3$, поэтому алгоритм перейдёт в конфигурацию $(\perp, *_3, i \& i, \perp, \varepsilon)$. Аналогично для магазинного символа $*_3$ и следующего входного символа i , алгоритм перейдёт в конфигурацию $(\perp, *_3, i_4, \& i, \perp, \varepsilon)$.

Шаг 3. Для строки таблицы, отмеченной символом $i_4, f(\&) = (C, 4)$, значит, необходимо выполнить свёртку по правилу $p_4: F \rightarrow i$.

Правая часть этого правила содержит только один символ, поэтому:

- Удаляем из магазина символ i_4 ;
- Определяем значение функции переходов для символа F из левой части правила p_4 в строке таблицы M , отмеченной символом $*_3$, который теперь стал верхним символом магазина. Значение функции переходов: $g(F) = F_5$.
- Записываем в выходную цепочку число 4 (номер примененного правила).

Алгоритм переходит в конфигурацию $(\perp, *_3, F_5, \& i, \perp, 4)$.

Продолжая, получим следующую последовательность тактов:

$(\perp, *_3, F_5, \& i, \perp, 4) \xrightarrow{} (\perp, *_3 L_3, \& i, \perp, 45) \xrightarrow{} (\perp, (F_1, \& i), \perp, 453) \xrightarrow{} (\perp, (F_1 \& i_1, i), \perp, 453) \xrightarrow{} (\perp, (F_1 \& i_4,), \perp, 453) \xrightarrow{} (\perp, (F_1 \& i_5,), \perp, 4534) \xrightarrow{} (\perp, (F_1 \& i_1 L_1,), \perp, 45345) \xrightarrow{} (\perp, (S_2,), \perp, 453451) \xrightarrow{} (\perp, (S_2)_2, \perp, 453451) \xrightarrow{} (\perp, S_2, \perp, 4534512) \xrightarrow{} \text{ДОПУСК}$.

Пусть анализируется цепочка $(*i \& & i)$, содержащая синтаксическую ошибку, последовательность тактов будет следующая:

$(\perp, (*i \& & i), \perp, \varepsilon) \xrightarrow{} (\perp, (*i \& & i), \perp, \varepsilon) \xrightarrow{} (\perp, (*_3, i \& i), \perp, \varepsilon) \xrightarrow{} (\perp, (*_3 i_4, \& i), \perp, \varepsilon) \xrightarrow{} (\perp, (*_3 F_5, \& i), \perp, 4) \xrightarrow{} (\perp, (*_3 L_3, \& i), \perp, 45) \xrightarrow{} (\perp, (F_1, \& i), \perp, 453) \xrightarrow{} (\perp, (F_1 \& i_1, i), \perp, 453) \xrightarrow{} \text{ОШИБКА}$.

Для строки таблицы, отмеченной символом $\&_1 f(\&) = \text{ПЕРЕНОС}$, но значение функции $g(\&) = \text{ОШИБКА}$. Алгоритм должен выдать сообщение об ошибке и завершить работу.

2. Построение LR(k) анализатора способом грамматических вхождений

Способ использования грамматических вхождений состоит из четырёх шагов:

Шаг 1. Определение грамматических вхождений (см. алгоритм фреймворка);

Шаг 2. Построение конечного автомата из грамматических вхождений;

Шаг 3. Построение управляющей таблицы.

Шаг 4. Применение алгоритма «перенос-свёртка».

Шаг 1. Определение грамматических вхождений

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

Пример

В приведённой выше грамматике с правилами:

$p_1: S \rightarrow F \& L,$
 $p_2: S \rightarrow (S),$
 $p_3: F \rightarrow * L,$
 $p_4: F \rightarrow i,$
 $p_5: L \rightarrow F$

можно, например, выделить грамматические вхождения $\&_{1,2}, \&_{2,3}, L_{3,2}$ и т.д.

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.

Шаг 2. Построение конечного автомата из грамматических вхождений

LR(0) ситуацией назовём запись $[A \rightarrow w_1 \bullet w_2]$, если $A \rightarrow w_1 w_2$ – правило КС-грамматики.

Пример.

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

$[S \rightarrow \bullet (S)]$
 $[S \rightarrow (\bullet S)]$
 $[S \rightarrow (S \bullet)]$
 $[S \rightarrow (S) \bullet]$

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

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 $A \rightarrow \alpha \cdot B\beta$ уже включена в $CLOSURE(I)$ и $B \rightarrow \gamma$ – правило грамматики, то добавить в множество $CLOSURE(I)$ ситуацию $B \rightarrow \cdot \gamma$ 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 $A \rightarrow \alpha \cdot B\beta$ 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 $B \rightarrow \cdot \gamma$.

Рассмотрим дополненную грамматику G_0 , содержащую следующие правила:

$$\begin{aligned} p_0: S' &\rightarrow S \\ p_1: S &\rightarrow F \& L, \\ p_2: S &\rightarrow (S), \\ p_3: F &\rightarrow * L, \\ p_4: F &\rightarrow i, \\ p_5: L &\rightarrow F \end{aligned}$$

Пусть вначале множество ситуаций $CLOSURE(I)$ содержит одну ситуацию $S' \rightarrow \cdot S$, т.е. $CLOSURE(I) = \{S' \rightarrow \cdot S\}$. Ситуация $S' \rightarrow \cdot S$ содержит точку перед символом S , поэтому по второму правилу в $CLOSURE(I)$ необходимо включить S -правила с точкой слева. Окончательно получим:

$$CLOSURE(I) = \{S' \rightarrow \cdot S, S \rightarrow \cdot F \& L, S \rightarrow \cdot (S), F \rightarrow \cdot * L, F \rightarrow \cdot i\}.$$

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 $[A \rightarrow \alpha X \beta]$, таких что $[A \rightarrow \alpha \cdot X \beta] \in I$

Пусть для грамматики G_0 некоторое множество $I = \{[S' \rightarrow S \cdot], [S \rightarrow F \cdot \& L]\}$.

Для вычисления значения $GOTO(I, \&)$ необходимо рассмотреть ситуации, в которых непосредственно сразу за точкой следует символ $\&$. Это ситуация $[S \rightarrow F \cdot \& L]$.

Перемещая точку за символ $\&$, построим множество ситуаций

$[S \rightarrow F \& \cdot L]$ и замыкание этого множества: $\{L \rightarrow \cdot F, F \rightarrow \cdot i, F \rightarrow \cdot * L\}$. Тогда $GOTO(I, \&) = \{[S \rightarrow F \& \cdot L], [L \rightarrow \cdot F], [F \rightarrow \cdot i], [F \rightarrow \cdot * L]\}$.

В [2] и [3] приведены алгоритмы, позволяющие построить каноническую систему множеств $LR(0)$ -ситуаций. 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, «не отмечено», то:

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

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

Пример

Построим систему множеств φ для нашей дополненной грамматики G_0 :

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

Определив множество $A_0 = CLOSURE(\{S' \rightarrow \cdot S\}) = \{S' \rightarrow \cdot S, S \rightarrow \cdot F \& L, S \rightarrow \cdot (S), F \rightarrow \cdot i, F \rightarrow \cdot * L\}$, включим его в систему φ в качестве «неотмеченного» множества.

A_0	$S' \rightarrow \cdot S$ $S \rightarrow \cdot F \& L$ $S \rightarrow \cdot (S)$ $F \rightarrow \cdot i$ $F \rightarrow \cdot * L$
-------	---

Теперь мы должны отметить множество A_0 и определить множества $GOTO(A_0, X)$ для всех символов X грамматики G_0 : $A_1 = GOTO(A_0, S) = CLOSURE(\{S' \rightarrow S \cdot\}) = \{S' \rightarrow S \cdot\}$

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

A_0	$S' \rightarrow S$ $S \rightarrow \bullet F \& L$ $S \rightarrow \bullet (S)$ $F \rightarrow \bullet i$ $F \rightarrow \bullet \bullet L$
A_1	$S' \rightarrow S \bullet$

Продолжая строить $GOTO(A_0, X)$, получаем:

$$A_2 = GOTO(A_0, ()) = CLOSURE([S \rightarrow (\bullet S)]) = \{S \rightarrow (\bullet S), S \rightarrow \bullet (S), S \rightarrow \bullet F \& L, F \rightarrow \bullet i, F \rightarrow \bullet \bullet L\}$$

$$A_3 = GOTO(A_0, i) = \{F \rightarrow i \bullet\}$$

$$A_4 = GOTO(A_0, *) = \{F \rightarrow \bullet \bullet L, L \rightarrow \bullet F, F \rightarrow \bullet \bullet L, F \rightarrow \bullet i\}$$

$$A_5 = GOTO(A_0, F) = \{S \rightarrow F \bullet \& L\}$$

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

A_0	$S' \rightarrow \bullet S$ $S \rightarrow \bullet F \& L$ $S \rightarrow \bullet (S)$ $F \rightarrow \bullet i$ $F \rightarrow \bullet \bullet L$
A_1	$S' \rightarrow S \bullet$
A_2	$S \rightarrow (\bullet S)$ $S \rightarrow \bullet (S)$ $S \rightarrow \bullet F \& L$ $F \rightarrow \bullet i$ $F \rightarrow \bullet \bullet L$
A_3	$F \rightarrow i \bullet$
A_4	$F \rightarrow \bullet \bullet L$ $L \rightarrow \bullet F$ $F \rightarrow \bullet \bullet L$ $F \rightarrow \bullet i$
A_5	$S \rightarrow F \bullet \& L$

Продолжая выполнять шаг 3 построения системы множеств, «отмечаем» множество A_1 . Для всех $X \in (V \cup \Sigma)$ множества $GOTO(A_1, X)$ пусты.

Теперь отмечаем множество A_2 . $A_6 = GOTO(A_2, S) = \{S \rightarrow (S \bullet)\}$

Добавляем множество A_6 в φ . $A_7 = GOTO(A_2, ()) = CLOSURE(S \rightarrow (\bullet S)) = \{S \rightarrow (\bullet S), S \rightarrow \bullet (S), S \rightarrow \bullet F \& L, F \rightarrow \bullet i, F \rightarrow \bullet \bullet L\}$

Такое множество (само множество A_2) уже есть в φ .

Множества $GOTO(A_2, F)$, $GOTO(A_2, i)$, $GOTO(A_2, *)$ также уже включены в φ . Для остальных символов $X \in (V \cup \Sigma)$ множества $GOTO(A_2, X)$ пусты.

Отмечаем множество A_3 . Для всех $X \in (V \cup \Sigma)$ множества $GOTO(A_3, X)$ пусты.

Продолжая аналогичным образом, включаем в систему φ множества:

$$A_7 = GOTO(A_4, L) = \{F \rightarrow \bullet \bullet L \bullet\}$$

$$A_8 = GOTO(A_4, F) = \{L \rightarrow F \bullet\}$$

$$A_9 = GOTO(A_5, \&) = \{S \rightarrow F \& \bullet L, L \rightarrow \bullet F, F \rightarrow \bullet i, F \rightarrow \bullet \bullet L\}$$

$$A_{10} = GOTO(A_6, *) = \{S \rightarrow (S \bullet)\}$$

$$A_{11} = GOTO(A_9, L) = \{S \rightarrow F \& L \bullet\}$$

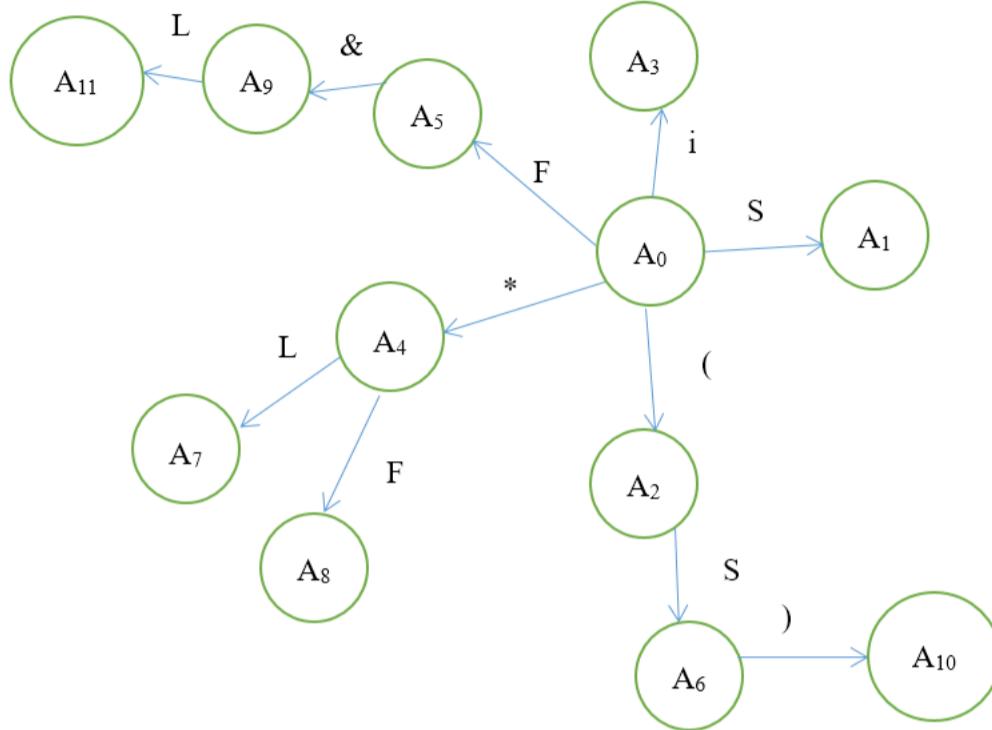
В результате получаем окончательную систему LR(0) – множеств φ :

A_0	$S' \rightarrow \bullet S$ $S \rightarrow \bullet F \& L$ $S \rightarrow \bullet (S)$ $F \rightarrow \bullet i$ $F \rightarrow \bullet \bullet L$
A_1	$S' \rightarrow S \bullet$
A_2	$S \rightarrow (\bullet S)$ $S \rightarrow \bullet (S)$ $S \rightarrow \bullet F \& L$ $F \rightarrow \bullet i$ $F \rightarrow \bullet \bullet L$
A_3	$F \rightarrow i \bullet$
A_4	$F \rightarrow \bullet \bullet L$ $L \rightarrow \bullet F$ $F \rightarrow \bullet \bullet L$ $F \rightarrow \bullet i$
A_5	$S \rightarrow F \bullet \& L$
A_6	$S \rightarrow (S \bullet)$
A_7	$F \rightarrow \bullet L \bullet$
A_8	$L \rightarrow F \bullet$
A_9	$S \rightarrow F \& \bullet L$ $L \rightarrow \bullet F$

	$F \rightarrow \bullet i$ $F \rightarrow \bullet * L$
A ₁₀	$S \rightarrow (S) \bullet$
A ₁₁	$S \rightarrow F \& L \bullet$

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

Диаграмма переходов ДКА для активных префиксов грамматики G0:



Переход в состояние, которому на диаграмме соответствует лист (вершина, не имеющая исходящих дуг) однозначно определяет функцию свёртки с возвратом к предыдущему состоянию.

Шаг 3. Построение управляющей таблицы.

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.

Управляющая таблица:

	()	&	I	*	\$	S	F	L
0	$\Pi(2)$			$\Pi(3)$	$\Pi(4)$			5	
1						допуск			
2							6		
3	C(4)	C(4)	C(4)	C(4)	C(4)	C(4)			
4								8	7
5			$\Pi(9)$						
6		$\Pi(10)$							
7	C(3)	C(3)	C(3)	C(3)	C(3)	C(3)			
8	C(5)	C(5)	C(5)	C(5)	C(5)	C(5)			
9			$\Pi(11)$						
10	C(2)	C(2)	C(2)	C(2)	C(2)				
11	C(1)	C(1)	C(1)	C(1)	C(1)				

Шаг 4. Применение алгоритма «перенос-свёртка».

Работа алгоритма описывается в терминах конфигураций, представляющих собой тройки вида $(\alpha T, ax, \pi)$, где αT – цепочка магазинных символов (T -верхний символ магазина), ax – необработанная часть входной цепочки, π – выход, построенный к настоящему моменту времени.

Рассмотрим последовательность тактов LR(k)-алгоритма при анализе входной цепочки abc\$. В магазине МП-автомата вместе с помещенным туда символами показаны и номера строк управляющей таблицы, в формате «символ»«номер».

Пример

Рассмотрим последовательность тактов алгоритма при анализе входной цепочки $i\&**i$

$(\perp, i\&**i$, \epsilon) \rightarrow (\perp 0i3, \&**i$, \epsilon) \rightarrow (\perp 0F5, \&**i$, 4) \rightarrow (\perp 0F5&9, **i$, 4) \rightarrow (\perp 0F5&9*4, *i$, 4) \rightarrow (\perp 0F5&9*4*4, i$, 4) \rightarrow (\perp 0F5&9*4*4i3, $, 4) \rightarrow (\perp 0F5&9*4*4F8, $, 44) \rightarrow (\perp 0F5&9*4*4L7, $, 445) \rightarrow (\perp 0F5&9*4F8, $, 4453) \rightarrow (\perp 0F5&9*4L7, $, 44535) \rightarrow (\perp 0F5&9F8, $, 445353) \rightarrow (\perp 0F5&9L11, $, 4453535) \rightarrow (\perp 0S1, $, 44535351) \rightarrow \text{ДОПУСК}$

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).

Принцип организации распознавателя входных цепочек языка, заданного грамматикой предшествования, основывается на том, что для каждой упорядоченной пары символов в грамматике устанавливается некоторое отношение, называемое отношением предшествования. В процессе разбора входной цепочки расширенный МП-автомат сравнивает текущий символ входной цепочки с одним из символов, находящихся на верхушке стека автомата. В процессе сравнения проверяется, какое из возможных отношений предшествования существует между этими двумя символами. В зависимости от найденного отношения выполняется либо сдвиг (перенос), либо свертка. При отсутствии отношения предшествования между символами алгоритм сигнализирует об ошибке.

Задача заключается в том, чтобы иметь возможность непротиворечивым образом определить отношения предшествования между символами грамматики. Если это возможно, то грамматика может быть отнесена к одному из классов грамматик предшествования.

Существует несколько видов грамматик предшествования. Они различаются по тому, какие отношения предшествования в них определены и между какими типами символов (терминальными или нетерминальными) могут быть установлены эти отношения. Кроме того, возможны незначительные модификации функционирования самого алгоритма «сдвиг-свертка» в распознавателях для таких грамматик (в основном на этапе выбора правила для выполнения свертки, когда возможны неоднозначности) [5, 6, 23, 65].

Выделяют следующие виды грамматик предшествования:

- простого предшествования;
- расширенного предшествования;
- слабого предшествования;
- смешанной стратегии предшествования;
- операторного предшествования.

Рассмотрим два наиболее простых и распространенных типа — грамматики простого и операторного предшествования.

Грамматики простого предшествования

Грамматикой простого предшествования называют такую приведенную (без циклов, бесплодных и недостижимых символов и λ -правил) КС-грамматику $G(VN, VT, P, S)$, $V = VT \cup VN$, в которой:

1. Для каждой упорядоченной пары терминальных и нетерминальных символов выполняется не более чем одно из трех отношений предшествования:

- $B_i =\bullet B_j (\forall B_i, B_j \in V)$, если и только если \exists правило $A \rightarrow xB_iB_jy \in P$, где $A \in VN$, $x, y \in V^*$;
- $B_i <\bullet B_j (\forall B_i, B_j \in V)$, если и только если \exists правило $A \rightarrow xB_iDy \in P$ и вывод $D \Rightarrow^* S_j z$, где $A, D \in VN$, $x, y, z \in V^*$;
- $B_i \bullet> B_j (\forall B_i, B_j \in V)$, если и только если \exists правило $A \rightarrow xCB_jy \in P$ и вывод $C \Rightarrow^* zB_i$ или \exists правило $A \rightarrow xCDy \in P$ и выводы $C \Rightarrow^* zB_i$ и $D \Rightarrow^* B_j w$, где $A, C, D \in VN$, $x, y, z, w \in V^*$.

2. Различные правила в грамматике имеют разные правые части (то есть в грамматике не должно быть двух различных правил с одной и той же правой частью).

Отношения $=\bullet$, $<\bullet$ и $\bullet>$ называют отношениями предшествования для символов. Отношение предшествования единственны для каждой упорядоченной пары символов. При этом между какими-либо двумя символами может и не быть отношения предшествования — это значит, что они не могут находиться рядом ни в одном элементе разбора синтаксически правильной цепочки. Отношения предшествования зависят от порядка, в котором стоят символы, и в этом смысле их нельзя путать со знаками математических операций — они не обладают ни свойством коммутативности, ни свойством ассоциативности. Например, если известно, что $B_i \bullet> B_j$, то не обязательно выполняется $B_j <\bullet B_i$ (поэтому знаки предшествования иногда помечают специальной точкой: $=\bullet$, $<\bullet$, $\bullet>$).

Для грамматик простого предшествования известны следующие полезные свойства:

- всякая грамматика простого предшествования является однозначной;
- легко проверить, является ли произвольная КС-грамматика грамматикой простого предшествования (для этого достаточно проверить рассмотренные выше свойства грамматик простого предшествования или воспользоваться алгоритмом построения матрицы предшествования, который рассмотрен далее).

Как и для многих других классов грамматик, для грамматик простого предшествования не существует алгоритма, который бы мог преобразовать произвольную КС-грамматику в грамматику простого предшествования (или доказать, что преобразование невозможно).

Метод предшествования основан на том факте, что отношения предшествования между двумя соседними символами распознаваемой строки соответствуют трем следующим вариантам:

- $B_i <\bullet B_{i+1}$, если символ B_{i+1} — крайний левый символ некоторой основы (это отношение между символами можно назвать «предшествует основе» или просто «предшествует»);
- $B_i \bullet> B_{i+1}$, если символ B_i — крайний правый символ некоторой основы (это отношение между символами можно назвать «следует за основой» или просто «следует»);

- $B_i =\bullet B_{i+1}$, если символы B_i и B_{i+1} принадлежат одной основе (это отношение между символами можно назвать «составляют основу»),

Исходя из этих соотношений, выполняется разбор строки для грамматики предшествования.

Суть принципа такого разбора можно пояснить на рис. 12.5. На нем изображена входная цепочка символов $\alpha\beta\gamma\delta$ в тот момент, когда выполняется свертка цепочки γ . Символ a является последним символом подцепочки α , а символ b — первым символом подцепочки β . Тогда, если в грамматике удастся установить непротиворечивые отношения предшествования, то в процессе выполнения разбора по алгоритму «сдвиг-свертка» можно всегда выполнять сдвиг до тех пор, пока между символом на верхушке стека и текущим символом входной цепочки существует отношение $<\bullet$ или $=\bullet$. А как только между этими символами будет обнаружено отношение $\bullet>$, так сразу надо выполнять свертку. Причем для выполнения свертки из стека надо выбирать все символы, связанные отношением $=\bullet$. То, что все различные правила в грамматике предшествования имеют различные правые части, гарантирует непротиворечивость выбора правила при выполнении свертки.

Таким образом, установление непротиворечивых отношений предшествования между символами грамматики в комплексе с несовпадающими правыми частями различных правил дает ответы на все вопросы, которые надо решить для организации работы алгоритма «сдвиг-свертка» без возвратов.

\perp_H	a	a	γ	b	β	δ	\perp_K
-----------	-----	-----	----------	-----	---------	----------	-----------

Рис. 12.5. Отношения между символами входной цепочки в грамматике предшествования

На основании отношений предшествования строят матрицу предшествования грамматики. Строки матрицы предшествования помечаются первыми (левыми) символами, столбцы — вторыми (правыми) символами отношений предшествования. В клетки матрицы на пересечении соответствующих столбца и строки помещаются знаки отношений. При этом пустые клетки матрицы говорят о том, что между данными символами нет ни одного отношения предшествования.

Матрицу предшествования грамматики сложно построить, опираясь непосредственно на определения отношений предшествования. Удобнее воспользоваться двумя дополнительными множествами — множеством крайних левых и множеством крайних правых символов относительно нетерминальных символов грамматики $G(VN, VT, P, S)$, $V = VT \cup VN$. Эти множества определяются следующим образом:

- $L(A) = \{X \mid \exists A \Rightarrow^* Xz\}, A \in VN, X \in V, z \in V^*$ — множество крайних левых символов относительно нетерминального символа A (цепочка z может быть и пустой цепочкой);
- $R(A) = \{X \mid \exists A \Rightarrow^* zX\}, A \in VN, X \in V, z \in V^*$ — множество крайних правых символов относительно нетерминального символа A .

Иными словами, множество крайних левых символов относительно нетерминального символа A — это множество всех крайних левых символов в цепочках, которые могут быть выведены из символа A . Аналогично, множество крайних правых символов относительно нетерминального символа A — это множество всех крайних правых символов в цепочках, которые могут быть выведены из символа A . Тогда отношения предшествования можно определить так:

- $B_i =\bullet B_j (\forall B_i, B_j \in V)$, если \exists правило $A \Rightarrow xB_iB_jy \in P$, где $A \in VN, x, y \in V^*$;
- $B_i <\bullet B_j (\forall B_i, B_j \in V)$, если \exists правило $A \Rightarrow xB_iDy \in P$ и $B_j \in L(D)$, где $A, D \in VN, x, y \in V^*$;
- $B_i \bullet> B_j (\forall B_i, B_j \in V)$, если \exists правило $A \Rightarrow xCB_jy \in P$ и $B_i \in R(C)$ или \exists правило $A \Rightarrow xCDy \in P$ и $B_i \in R(C), B_j \in L(D)$, где $A, C, D \in VN, x, y \in V^*$.

Такое определение отношений удобнее на практике, так как не требует построения выводов, а множества $L(A)$ и $R(A)$ могут быть построены для каждого нетерминального символа $A \in VN$ грамматики $G(VN, VT, P, S)$, $V = VT \cup VN$ по очень простому алгоритму.

Шаг 1. $\forall A \in VN$:

$$R_0(A) = \{X \mid A \Rightarrow yX, X \in V, y \in V^*\}, L_0(A) = \{X \mid A \Rightarrow Xy, X \in V, y \in V^*\}, i := 1.$$

Для каждого нетерминального символа A ищем все правила, содержащие A в левой части. Во множество $L(A)$ включаем самый левый символ из правой части правил, а во множество $R(A)$ — самый крайний правый символ из правой части. Переходим к шагу 2.

Шаг 2. $\forall A \in VN$:

$$R_i(A) = R_{i-1}(A) \cup R_{i-1}(B), \forall B \in (R_{i-1}(A) \cap VN),$$

$$L_i(A) = L_{i-1}(A) \cup L_{i-1}(B), \forall B \in (L_{i-1}(A) \cap VN).$$

Для каждого нетерминального символа A : если множество $L(A)$ содержит нетерминальные символы грамматики A', A'', \dots , то его надо дополнить символами, входящими в соответствующие множества $L(A'), L(A''), \dots$ и не входящими в $L(A)$. Ту же операцию надо выполнить для $R(A)$.

Шаг 3. Если $\exists A \in VN: R_i(A) \neq R_{i-1}(A)$ или $L_i(A) \neq L_{i-1}(A)$, то $i := i + 1$ и вернуться к шагу 2, иначе построение закончено: $R(A) = R_i(A)$ и $L(A) = L_i(A)$.

Если на предыдущем шаге хотя бы одно множество $L(A)$ или $R(A)$ для некоторого символа грамматики изменилось, то надо вернуться к шагу 2, иначе построение закончено.

После построения множеств $L(A)$ и $R(A)$ по правилам грамматики создается матрица предшествования. Матрицу предшествования дополняют символами \perp_H и \perp_K (начало и конец цепочки). Для них определены следующие отношения предшествования:

$\perp_H <\bullet X, \forall a \in V$, если $\exists S \Rightarrow^* Xy, \text{ где } S \in VN, y \in V^*$ или (с другой стороны) если $X \in L(S)$;

$\perp_K \bullet> X, \forall a \in V$, если $\exists S \Rightarrow^* yX, \text{ где } S \in VN, y \in V^*$ или (с другой стороны) если $X \in R(S)$;

Здесь S — целевой символ грамматики.

Матрица предшествования служит основой для работы распознавателя языка, заданного грамматикой простого предшествования.

Алгоритм «сдвиг-свертка» для грамматики простого предшествования

Данный алгоритм выполняется расширенным МП-автоматом с одним состоянием. Отношения предшествования служат для того, чтобы определить в процессе выполнения алгоритма, какое действие — сдвиг или свертка — должно выполняться на каждом шаге алгоритма, и однозначно выбрать цепочку для свертки. Однозначный выбор правила при свертке обеспечивается за счет различия правых частей всех правил грамматики. В начальном состоянии автоматачитывающая головка обозревает первый символ входной цепочки, в стеке МП-автомата находится символ \perp_h (начало цепочки), в конец цепочки помещен символ \perp_k (конец цепочки). Символы \perp_h и \perp_k введены для удобства работы алгоритма, в язык, заданный исходной грамматикой, они не входят.

Разбор считается законченным (алгоритм завершается), есличитывающая головка автомата обозревает символ \perp_k и при этом больше не может быть выполнена свертка. Решение о принятии цепочки зависит от содержимого стека. Автомат принимает цепочку, если в результате завершения алгоритма в стеке находятся начальный символ грамматики S и символ \perp_h . Выполнение алгоритма может быть прервано, если на одном из его шагов возникнет ошибка. Тогда входная цепочка не принимается. Алгоритм состоит из следующих шагов.

Шаг 1. Поместить в верхушку стека символ \perp_h , считающую головку — в начало входной цепочки символов.

Шаг 2. Сравнить с помощью отношения предшествования символ, находящийся на вершине стека (левый символ отношения), с текущим символом входной цепочки, обозреваемымчитывающей головкой (правый символ отношения).

Шаг 3. Если имеет место отношение $<\bullet$ или $=\bullet$, то произвести сдвиг (перенос текущего символа из входной цепочки в стек и сдвигчитывающей головки на один шаг вправо) и вернуться к шагу 2. Иначе перейти к шагу 4.

Шаг 4. Если имеет место отношение $\bullet>$, то произвести свертку. Для этого надо найти на вершине стека все символы, связанные отношением $=\bullet$ («основу»), удалить эти символы из стека. Затем выбрать из грамматики правило, имеющее правую часть, совпадающую с основой, и поместить в стек левую часть выбранного правила (если символов, связанных отношением $=\bullet<\$]command>$, на верхушке стека нет, то в качестве основы используется один, самый верхний символ стека). Если правило, совпадающее с основой, найти не удалось, то необходимо прервать выполнение алгоритма и сообщить об ошибке, иначе, если разбор не закончен, то вернуться к шагу 2.

Шаг 5. Если не установлено ни одно отношение предшествования между текущим символом входной цепочки и символом на верхушке стека, то надо прервать выполнение алгоритма и сообщить об ошибке.

Ошибка в процессе выполнения алгоритма возникает, когда невозможно выполнить очередной шаг — например, если не установлено отношение предшествования между двумя сравниваемыми символами (на шагах 2 и 4) или если не удается найти нужное правило в грамматике (на шаге 4). Тогда выполнение алгоритма немедленно прерывается.

Пример 1. <https://studfiles.net/preview/4599549/page:23/>

https://life-prog.ru/view_psp.php?id=13

Дана грамматика $G(\{a, (,)\}, \{S, R\}, P, S)$, с правилами P :

1) $S \rightarrow (R|a); 2) R \rightarrow Sa$. Построить распознаватель для строки $((aa)a)a\perp_k$.

Шаг 1. Построим множества крайних левых и крайних правых символов $L(A)$ и $R(A)$ относительно всех нетерминальных символов грамматики (таблица 1.1).

Шаг	$L_i(A)$	$R_i(A)$
0	$L_0(S)=\{(, a\}$ $L_0(R)=\{S\}$	$R_0(S)=\{R, a\}$ $R_0(R)=\{\})\}$
1	$L_1(S)=\{(, a\}$ $L_1(R)=\{S, (, a\}$	$R_1(S)=\{R, a,)\}$ $R_1(R)=\{\})\}$
2	$L_2(S)=\{(, a\}$ $L_2(R)=\{S, (, a\}$	$R_2(S)=\{R, a,)\}$ $R_2(R)=\{\})\}$

Резул ьтат	$L(S)=\{(), a\}$	$R(S)=\{R, a,)\}$
	$L(R)=\{S, (), a\}$	$R(R)=\{\}\}$

Таблица 1.1 – Построение множеств $L(A)$ и $R(A)$ для грамматики G

Шаг 2. На основе построенных множеств и правил вывода грамматики составим матрицу предшествования символов (таблица 1.2).

Символы	S	R	a	()	\perp_k
S			=•			
R			•>			•>
a			•>		=•	•>
(<•	=•	<•	<•		
)			•>			•>
\perp_h			<•	<•		

Таблица 1.2 – Матрица предшествования символов грамматики

В правиле грамматики $S \rightarrow (R$ символ $($ стоит слева от нетерминального символа R . Во множество $L(R)$ входят символы $S, (, a$. Ставим знак $<•$ в клетках матрицы, соответствующих этим символам, в строке для символа $($.

В правиле грамматики $R \rightarrow Sa$ символ a стоит справа от нетерминального символа S . Во множество $R(S)$ входят символы $R, a,)$. Ставим знак $•>$ в клетках матрицы, соответствующих этим символам, в столбце для символа a .

В строке символа \perp_h ставим знак $<•$ в клетках символов, входящих во множество $L(S)$, т.е. символов $(, a$. В столбце символа \perp_k ставим знак $•>$ в клетках, входящих во множество $R(S)$, т.е. символов $R, a,)$.

В клетках, соответствующих строке символа S и столбцу символа a , ставим знак $=$; т.к. существует правило $R \rightarrow Sa$, в котором эти символы стоят рядом. По тем же соображениям ставим знак $=$ в клетках строки a и столбца $)$, а также строки $($ и столбца R .

Шаг 3. Функционирование распознавателя для цепочки $((aa)a)a$ показано в таблице 1.3.

Шаг	Стек	Входной буфер	Действие
1	\perp_h	$((aa)a)a\perp_k$	сдвиг
2	$\perp_h($	$((aa)a)a\perp_k$	сдвиг
3	$\perp_h(($	$(aa)a)a\perp_k$	сдвиг
4	$\perp_h((($	$aa)a)a\perp_k$	сдвиг
5	$\perp_h(((a$	$a)a)a\perp_k$	свертка $S \rightarrow a$
6	$\perp_h(((S$	$a)a)a\perp_k$	сдвиг
7	$\perp_h(((Sa$	$)a)a\perp_k$	сдвиг
8	$\perp_h(((Sa)$	$a)a\perp_k$	свертка $R \rightarrow Sa$
9	$\perp_h(((R$	$a)a\perp_k$	свертка $S \rightarrow (R$
10	$\perp_h((S$	$a)a\perp_k$	сдвиг
11	$\perp_h((Sa$	$)a\perp_k$	сдвиг
12	$\perp_h((Sa)$	$a\perp_k$	свертка $R \rightarrow Sa$
13	$\perp_h((R$	$a\perp_k$	свертка $S \rightarrow (R$
14	$\perp_h(S$	$a\perp_k$	сдвиг
15	$\perp_h(Sa$	$)\perp_k$	сдвиг
16	$\perp_h(Sa)$	\perp_k	свертка $R \rightarrow Sa$
17	$\perp_h(R$	\perp_k	свертка $S \rightarrow (R$
18	$\perp_h S$	\perp_k	строка принята

Таблица 1.3. – Алгоритм работы распознавателя цепочки $((aa)a)a$

Шаг 4. Получили следующую цепочку вывода:

$$S \Rightarrow (R \Rightarrow (Sa) \Rightarrow ((Ra) \Rightarrow ((Sa)a) \Rightarrow (((Ra)a) \Rightarrow (((Sa)a)a) \Rightarrow (((aa)a)a).$$

Пример 2. $I \rightarrow TR \mid T$

$$R \rightarrow +T \mid -T \mid +TR \mid -TR$$

$$T \rightarrow a \mid b$$

Построим крайние левые и правые:

1. $L(I) = \{ T \}$
2. $L(R) = \{ +, - \}$
3. $L(T) = \{ a, b \}$
4. $R(I) = \{ R, T \}$
5. $R(R) = \{ T, R \}$
6. $R(T) = \{ a, b \}$

Дополняем к Т входящие a,b:

7. $L(I) = \{ T, a, b \}$
8. $L(R) = \{ +, - \}$
9. $L(T) = \{ a, b \}$
10. $R(I) = \{ R, T, a, b \}$
11. $R(R) = \{ T, R, a, b \}$
12. $R(T) = \{ a, b \}$

	T	R	+	-	a	b	Ik
T		=•	<•	<•			•>
R							•>
+	=•				<•	<•	
-	=•				<•	<•	
a			•>	•>	•>		•>
b			•>	•>	•>		•>
In	<•				<•	<•	

Процесс распознавания работает как LR анализатор:

a+b-a	In	Π
a+b-a	Ina	C
+b-a	InT	Π
b-a	InT+	Π
-a	InT+b	C
-a	InT+T	Π
a	InT+T-	Π
Ik	InT+T-a	C
Ik	InT+T-T	C
Ik	InT+TR	C
Ik	InTR	C
Ik	InI	C

Пример 3.

$$\begin{aligned} I &\rightarrow TS; \\ T &\rightarrow \text{int} \mid \text{char} \\ S &\rightarrow *BR \mid BR \mid B \\ R &\rightarrow , \mid S \\ B &\rightarrow i \mid j \end{aligned}$$

Построим крайние левые и правые:

$$\begin{aligned} L(I) &= \{ T, \text{int}, \text{char} \} \\ L(T) &= \{ \text{int}, \text{char} \} \\ L(S) &= \{ *, B, i, j \} \end{aligned}$$

$L(R) = \{,\}$
 $L(B) = \{i,j\}$
 $R(I) = \{\}$
 $R(T) = \{int, char\}$
 $R(S) = \{R, B, i, j, S\}$
 $R(R) = \{S, R, B, i, j\}$
 $R(B) = \{I, j\}$

Строим таблицы:

	T	S	;	int	char	*	B	R	,	i	j	Ik
T	=•				<•	<•			<•	<•		
S		=•			•>							
;												
Int	•>				•>	•>			•>	•>		
char	•>				•>	•>			•>	•>		
*						=•			<•	<•		
B		•>					=•	<•				
R		•>										
,												
I		•>						•>	•>			
J		•>						•>	•>			
In	<•			<•	<•							

В таблице имеет два предшествования. Это на пересечении ; и S.

Следовательно правило нужно переписать:

$I \rightarrow TS;$
 $T \rightarrow int \mid char$
 $S \rightarrow , *BR \mid , BR \mid , B$
 $B \rightarrow i \mid j$

Тогда буква S должна исчезнуть.

$R(R) = \{R, B, i, j\}$

Грамматики операторного предшествования

Операторной грамматикой называется КС-грамматика без λ -правил, в которой правые части всех правил не содержат смежных нетерминальных символов. Для операторной грамматики отношения предшествования можно задать на множестве терминальных символов (включая символы \perp_H и \perp_K).

Грамматикой операторного предшествования называется операторная КС-грамматика $G(VN, VT, P, S)$, $V = VT \cup VN$, для которой выполняются следующие условия:

1. Для каждой упорядоченной пары терминальных символов выполняется не более чем одно из трех отношений предшествования:

- $a =• b$, если и только если существует правило $A \rightarrow xaby \in P$ или правило $A \rightarrow xaCby$, где $a, b \in VT$, $A, C \in VN$, $x, y \in V^*$;
- $a <• b$, если и только если существует правило $A \rightarrow xaCy \in P$ и вывод $C \Rightarrow^* bz$ или вывод $C \Rightarrow^* Dbz$, где $a, b \in VT$, $A, C, D \in VN$, $x, y, z \in V^*$;
- $a •> b$, если и только если существует правило $A \rightarrow xCBy \in P$ и вывод $C \Rightarrow^* za$ или вывод $C \Rightarrow^* zaD$, где $a, b \in VT$, $A, C, D \in VN$, $x, y, z \in V^*$.

2. Различные порождающие правила имеют разные правые части, λ -правила отсутствуют.

Отношения предшествования для грамматик операторного предшествования определены таким образом, что для них выполняется еще одна особенность — правила грамматики операторного предшествования не могут содержать двух смежных нетерминальных символов в правой части. То есть в грамматике операторного предшествования $G(VN, VT, P, S)$, $V = VT \cup VN$ не может быть ни одного правила вида: $A \rightarrow xBCy$, где $A, B, C \in VN$, $x, y \in V^*$ (здесь x и y — это произвольные цепочки символов, могут быть и пустыми).

Для грамматик операторного предшествования также известны следующие свойства:

- всякая грамматика операторного предшествования задает детерминированный КС-язык (но не всякая грамматика операторного предшествования при этом является однозначной!);
- легко проверить, является ли произвольная КС-грамматика грамматикой операторного предшествования (точно так же, как и для простого предшествования).

Как и для многих других классов грамматик, для грамматик операторного предшествования не существует алгоритма, который бы мог преобразовать произвольную КС-грамматику в грамматику операторного предшествования (или доказать, что преобразование невозможно).

Принцип работы распознавателя для грамматики операторного предшествования аналогичен грамматике простого предшествования, но отношения предшествования проверяются в процессе разбора только между терминальными символами.

Для грамматики данного вида на основе установленных отношений предшествования также строится матрица предшествования, но она содержит только терминальные символы грамматики.

Для построения этой матрицы удобно ввести множества крайних левых и крайних правых терминалных символов относительно нетерминального символа A - $L^t(A)$ или $R^t(A)$:

- $L^t(A) = \{t \mid \exists A \Rightarrow^* tz \text{ или } \exists A \Rightarrow^* Ctz\}$, где $t \in VT$, $A, C \in VN$, $z \in V^*$;
- $R^t(A) = \{t \mid \exists A \Rightarrow^* zt \text{ или } \exists A \Rightarrow^* ztC\}$, где $t \in VT$, $A, C \in VN$, $z \in V^*$.

Тогда определения отношений операторного предшествования будут выглядеть так:

- $a = \bullet$ Ь, если \exists правило $A \rightarrow xaby \in P$ или правило $U \rightarrow xaCby$, где $a, b \in VT$, $A, C \in VN$, $x, y \in V^*$;
- $a < \bullet$ Ь, если \exists правило $A \rightarrow xCy \in P$ и $b \in L^t(C)$, где $a, b \in VT$, $A, C \in VN$, $x, y \in V^*$;
- $a \bullet > \bullet$ Ь, если \exists правило $A \rightarrow xCby \in P$ и $a \in R^t(C)$, где $a, b \in VT$, $A, C \in VN$, $x, y \in V^*$.

В данных определениях цепочки символов x, y, z могут быть и пустыми цепочками.

Для нахождения множеств $L^t(A)$ и $R^t(A)$ предварительно необходимо выполнить построение множеств $L(A)$ и $R(A)$, как это было рассмотрено ранее. Далее для построения $L^t(A)$ и $R^t(A)$ используется следующий алгоритм:

Шаг 1. $\forall A \in VN$:

$$R_0^t(A) = \{t \mid A \rightarrow ytB \text{ или } A \rightarrow yt, t \in VT, B \in VN, y \in V^*\},$$
$$L_0^t(A) = \{t \mid A \rightarrow Bty \text{ или } A \rightarrow ty, t \in VT, B \in VN, y \in V^*\}.$$

Для каждого нетерминального символа A ищем все правила, содержащие A в левой части. Во множество $L(A)$ включаем самый левый терминальный символ из правой части правил, игнорируя нетерминальные символы, а во множество $R(A)$ — самый крайний правый терминальный символ из правой части правил. Переходим к шагу 2.

Шаг 2. $\forall A \in VN$:

$$R_i^t(A) = R_{i-1}^t(A) \cup R_{i-1}^t(B), \forall B \in (R(A) \cap VN),$$
$$L_i^t(A) = L_{i-1}^t(A) \cup L_{i-1}^t(B), \forall B \in (L(A) \cap VN).$$

Для каждого нетерминального символа A: если множество $L(A)$ содержит нетерминальные символы грамматики A' , A'' , ..., то его надо дополнить символами, входящими в соответствующие множества $L^t(A')$, $L^t(A'')$, ... и не входящими в $L^t(A)$. Ту же операцию надо выполнить для множеств $R(A)$ и $R^t(A)$.

Шаг 3. Если $\exists A \in VN$: $R_i^t(A) \neq R_{i-1}^t(A)$ или $L_i^t(A) \neq L_{i-1}^t(A)$, то $i := i+1$ и вернуться к шагу 2, иначе построение закончено: $R(A) = R_i^t(A)$ и $L(A) = L_i^t(A)$.

Если на предыдущем шаге хотя бы одно множество $L^t(A)$ или $R^t(A)$ для некоторого символа грамматики изменилось, то надо вернуться к шагу 2, иначе построение закончено.

Для практического использования матрицу предшествования дополняют символами \perp_H и \perp_K (начало и конец цепочки). Для них определены следующие отношения предшествования:

$\perp_H < \bullet a$, $\forall a \in VT$, если $\exists S \Rightarrow^* ax$ или $\exists S \Rightarrow^* Cax$, где $S, C \in VN$, $x \in V^*$ или если $a \in L^t(S)$;

$\perp_K \bullet > a$, $\forall a \in VT$, если $\exists S \Rightarrow^* xa$ или $\exists S \Rightarrow^* xaC$, где $S, C \in VN$, $x \in V^*$ или если $a \in R^t(S)$.

Здесь S — целевой символ грамматики.

Матрица предшествования служит основой для работы распознавателя языка, заданного грамматикой операторного предшествования. Поскольку она содержит только терминальные символы, то, следовательно, будет иметь меньший размер, чем аналогичная матрица для грамматики простого предшествования. Следует отметить, что напрямую сравнивать матрицы двух грамматик нельзя — не всякая грамматика простого предшествования является грамматикой операторного предшествования, и наоборот. Например, рассмотренная далее в примере грамматика операторного предшествования не является грамматикой простого предшествования (читатели могут это проверить самостоятельно). Однако если две грамматики эквивалентны и задают один и тот же язык, то их множества терминалных символов должны совпадать. Тогда можно утверждать, что размер матрицы для грамматики операторного предшествования всегда будет меньше, чем размер матрицы эквивалентной ей грамматики простого предшествования.

Все, что было сказано выше о способах хранения матриц для грамматик простого предшествования, в равной степени относится также и к грамматикам операторного предшествования, с той только разницей, что объем хранимой матрицы будет меньше.

Алгоритм «сдвиг-свертка» для грамматики операторного предшествования

Алгоритм выполняется расширенным МП-автоматом и имеет те же условия завершения и обнаружения ошибок. Основное отличие состоит в том, что при определении отношения предшествования этот алгоритм не принимает во внимание находящиеся в стеке нетерминальные символы и при сравнении ищет ближайший к верхушке стека терминальный символ. Однако после выполнения сравнения и определения границ основы при поиске правила в грамматике нетерминальные символы следует, безусловно, принимать во внимание.

Алгоритм состоит из следующих шагов.

Шаг 1. Поместить в верхушку стека символ \perp_H , считающую головку — в начало входной цепочки символов.

Шаг 2. Сравнить с помощью отношения предшествования терминальный символ, ближайший к вершине стека (левый символ отношения), с текущим символом входной цепочки, обозреваемым считающей головкой (правый символ отношения). При этом из стека надо выбрать самый верхний терминальный символ, игнорируя все возможные нетерминальные символы.

Шаг 3. Если имеет место отношение $< \bullet$ или $= \bullet$, то произвести сдвиг (перенос текущего символа из входной цепочки в стек и сдвиг считающей головки на один шаг вправо) и вернуться к шагу 2. Иначе перейти к шагу 4.

Шаг 4. Если имеет место отношение $\bullet >$, то произвести свертку. Для этого надо найти на вершине стека все терминальные символы, связанные отношением $= \bullet$ («основу»), а также все соседствующие с ними нетерминальные символы (при определении отношения нетерминальные символы игнорируются). Если терминальных символов, связанных отношением $= \bullet$, на верхушке стека нет, то в качестве основы используется один, самый верхний в стеке терминальный символ стека. Все (и терминальные, и нетерминальные) символы, составляющие основу, надо удалить из

стека, а затем выбрать из грамматики правило, имеющее правую часть, совпадающую с основой, и поместить в стек левую часть выбранного правила. Если правило, совпадающее с основой, найти не удалось, то необходимо прервать выполнение алгоритма и сообщить об ошибке, иначе, если разбор не закончен, то вернуться к шагу 2.

Шаг 5. Если не установлено ни одно отношение предшествования между текущим символом входной цепочки и самым верхним терминальным символом в стеке, то надо прервать выполнение алгоритма и сообщить об ошибке.

Конечная конфигурация данного МП-автомата совпадает с конфигурацией при распознавании цепочек грамматик простого предшествования.

Пример

Дано: $G(\{(), \wedge, \&, \sim, a\}, \{S, T, E, F\}, P, S)$, где

$$P: S \rightarrow S^*T \mid T$$

$$T \rightarrow T \& E \mid E$$

$$E \rightarrow \sim E \mid F$$

$$F \rightarrow (E) \mid a$$

Построить: распознаватель для G .

Построение множеств $L(A)$ и $R(A)$:

i	$L_i(A)$	$R_i(A)$
0	$L_0(S)=\{S, T\}$ $L_0(T)=\{T, E\}$ $L_0(E)=\{\sim, F\}$ $L_0(F)=\{(), a\}$	$R_0(S)=\{T\}$ $R_0(T)=\{E\}$ $R_0(E)=\{E, F\}$ $R_0(F)=\{(), a\}$
1	$L_1(S)=\{S, T, E\}$ $L_1(T)=\{T, E, \sim, F\}$ $L_1(E)=\{\sim, F, (), a\}$ $L_1(F)=\{(), a\}$	$R_1(S)=\{T, E\}$ $R_1(T)=\{E, F\}$ $R_1(E)=\{E, F, (), a\}$ $R_1(F)=\{(), a\}$
2	$L_2(S)=\{S, T, E, \sim, F, (), a\}$ $L_2(T)=\{T, E, \sim, F, (), a\}$ $L_2(E)=\{\sim, F, (), a\}$ $L_2(F)=\{(), a\}$	$R_2(S)=\{T, E, F, (), a\}$ $R_2(T)=\{E, F, (), a\}$ $R_2(E)=\{E, F, (), a\}$ $R_2(F)=\{(), a\}$
3	$L_3(S)=\{S, T, E, \sim, F, (), a\}$ $L_3(T)=\{T, E, \sim, F, (), a\}$ $L_3(E)=\{\sim, F, (), a\}$ $L_3(F)=\{(), a\}$	$R_3(S)=\{T, E, F, (), a\}$ $R_3(T)=\{E, F, (), a\}$ $R_3(E)=\{E, F, (), a\}$ $R_3(F)=\{(), a\}$

Построение множеств $L'(A)$ и $R'(A)$:

i	$L'(A)$	$R'(A)$
0	$L'_0(S)=\{\wedge\}$ $L'_0(T)=\{\&\}$ $L'_0(E)=\{\sim\}$ $L'_0(F)=\{(), a\}$	$R'_0(S)=\{\wedge\}$ $R'_0(T)=\{\&\}$ $R'_0(E)=\{\sim\}$ $R'_0(F)=\{(), a\}$
1	$L'_1(S)=\{\wedge, \&, \sim, (), a\}$ $L'_1(T)=\{\&, \sim, (), a\}$ $L'_1(E)=\{\sim, (), a\}$ $L'_1(F)=\{(), a\}$	$R'_1(S)=\{\wedge, \&, \sim, (), a\}$ $R'_1(T)=\{\&, \sim, (), a\}$ $R'_1(E)=\{\sim, (), a\}$ $R'_1(F)=\{(), a\}$
2	$L'_2(S)=\{\wedge, \&, \sim, (), a\}$ $L'_2(T)=\{\&, \sim, (), a\}$ $L'_2(E)=\{\sim, (), a\}$ $L'_2(F)=\{(), a\}$	$R'_2(S)=\{\wedge, \&, \sim, (), a\}$ $R'_2(T)=\{\&, \sim, (), a\}$ $R'_2(E)=\{\sim, (), a\}$ $R'_2(F)=\{(), a\}$

Матрица операторного предшествования символов грамматики :

Символы	\wedge	$\&$	\sim	$($	$)$	a	$\perp_{\hat{e}}$
\wedge	$\bullet >$	$< \bullet$	$< \bullet$	$< \bullet$	$\bullet >$	$< \bullet$	$\bullet >$
$\&$	$\bullet >$	$\bullet >$	$< \bullet$	$< \bullet$	$\bullet >$	$< \bullet$	$\bullet >$
\sim	$\bullet >$	$\bullet >$	$< \bullet$	$< \bullet$	$\bullet >$	$< \bullet$	$\bullet >$
$($	$< \bullet$	$< \bullet$	$< \bullet$	$< \bullet$	$= \bullet$	$< \bullet$	
$)$	$\bullet >$	$\bullet >$			$\bullet >$		$\bullet >$
a	$\bullet >$	$\bullet >$			$\bullet >$		$\bullet >$
$\perp_{\hat{e}}$	$< \bullet$	$< \bullet$	$< \bullet$	$< \bullet$		$< \bullet$	

Для \wedge находящейся в правиле вывода слева от нетерминала T , во множество $L'(T)$ входят символы $\&, \sim, (), a$, значит в строке матрицы для \wedge ставим знак меньшего предшествования в позициях этих символов. С другой стороны этот символ \wedge находится справа от S . Во множество $R'(S)$ входят символы $\wedge, \&, \sim, (), a$, значит знак большего предшествования

ставится в столбце для \wedge в позициях этих символов. Символы (и) в правиле вывода находятся рядом, поэтому в позиции этих символов ставится знак равного предшествования (игнорируя нетерминал E).

++++++ Было

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 пару 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, ..., ОШИБКА}. Если $g(a, w) = i$, то правая часть i -го правила является суффиксом цепочки a .

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 $f(\alpha, aw) = \text{ПЕРЕНОС}$ (Π), 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 $f(\alpha\beta, w) = \text{СВЕРТКА}$ (C), $g(\alpha\beta, w) = i$ 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 $(\alpha, w, \pi) \vdash^S \text{ДОПУСК}$ (D).

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 $(\perp, w, \varepsilon) \vdash^* (\perp S, \varepsilon, \pi) \vdash^S \text{ДОПУСК}$, 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 правовыводимой 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	Π	Π			
αB	C	C			
Aa	C	C		Π	
Ab	Π		Π		C
Ac	Π		Π		
Ad	C	C			
\perp		Π			

^S | | | | | Д

Функция свертки

g	ax	bx	x	ε
$\perp bAb$				1
$\alpha A ad$		4		
αcB	2	2		
Aa			3	

2. Разберем с помощью построенного алгоритма λ входную цепочку $bcaadb$.

Начальная конфигурация алгоритма (\perp , bcaadb, ε).

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: $(\perp bca, adb, 3) \vdash^r (\perp bca, db, 3) \vdash^r (\perp bca, db, 3) \vdash^s (\perp bca, db, 34)$

$\vdash^s(\perp bA, b, 342) \vdash^r(\perp bAb, \varepsilon, 342) \vdash^r(\perp S, \varepsilon, 3421) \vdash^s$ 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 накладываются 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 $q_i a_i \rightarrow q_k \{ B, L, R \}$,

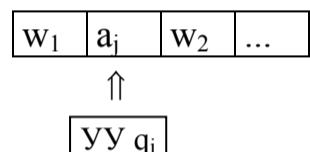
$\delta = Q \times \Sigma \times Q \times (\Sigma \cup \{ B, L, R \})$ $q_i a_j \rightarrow q_k \{ B, L, R \}$

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 $w_1, w_2 \in \Sigma^*$, $a_j \in \Sigma$, $q_i \in Q$.

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, управляющему устройству (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_j 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 $L = \{w \in \Sigma^* \mid q_0 \text{ и } w \Rightarrow^* C_f\}$, где C_f - 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 $\alpha = \xi_1 A \xi_2$, $\beta = \xi_1 \gamma \xi_2$, $A \subset V$; $\gamma \subset (T \cup V)^+$; $\xi_1, \xi_2 \subset (T \cup V)^*$.

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.

5. Соотношения между классами грамматик и языками

Отношения между классами КС-грамматик

В общем случае можно выделить правоанализируемые и левоанализируемые КС-грамматики. Об этих двух принципиально разных классах грамматик уже говорилось выше: первые предполагают построение левостороннего (восходящего) распознавателя, вторые — правостороннего (нисходящего). Это вовсе не значит, что для КС-языка, заданного, например, некоторой левоанализируемой грамматикой, невозможно построить расширенный МП-автомат, который порождает правосторонний вывод. Указанное разделение грамматик относится только к построению на их основе детерминированных МП-автоматов и детерминированных расширенных МП-автоматов. Только эти типы автоматов представляют интерес при создании компиляторов и анализе входных цепочек языков программирования. Недетерминированные автоматы, порождающие как левосторонние, так и правосторонние выводы, можно построить в любом случае для языка заданного любой КС-грамматикой, но для создания компилятора такие автоматы интереса не представляют (см. раздел «Распознаватели КС-языков. Автоматы с магазинной памятью»).

На рис. 12.10 изображена условная схема, дающая представление о соотношении классов левоанализируемых и правоанализируемых КС-грамматик [5, 6, т. 2, 42, 65].

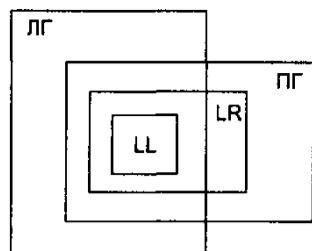


Рис. 12.10. Соотношение классов левоанализируемых и правоанализируемых КС-грамматик

Интересно, что классы левоанализируемых и правоанализируемых грамматик являются несопоставимыми. То есть существуют левоанализируемые КС-грамматики, на основе которых нельзя построить детерминированный расширенный МП-автомат, порождающий правосторонний вывод; и наоборот — существуют правоанализируемые КС-грамматики, не допускающие построение МП-автомата, порождающего левосторонний вывод. Конечно, существуют грамматики, подпадающие под оба класса и допускающие построение детерминированных автоматов как с правосторонним, так и с левосторонним выводом.

Следует помнить также, что все упомянутые классы КС-грамматик — это счетные, но бесконечные множества. Нельзя построить и рассмотреть все возможные левоанализируемые грамматики или даже все возможные LL(1)-грамматики. Сопоставление классов КС-грамматик производится исключительно на основе анализа структуры их правил. Только на основании такого рода анализа произвольная КС-грамматика может быть отнесена в тот или иной класс (или несколько классов).

Все это тем более интересно, если вспомнить, что рассмотренный в данном курсе класс левоанализируемых LL-грамматик является собственным подмножеством класса LR-грамматик: любая LL-грамматика является LR-грамматикой, но не наоборот — существуют LR-грамматики, которые не являются LL-грамматиками. Этот факт также нашел свое отражение в схеме на рис. 12.10. Значит, любая LL-грамматика является правоанализируемой, но существуют также и другие левоанализируемые грамматики, не попадающие в класс правоанализируемых грамматик.

Для LL(k)-грамматик, составляющих класс LL-грамматик, интересна еще одна особенность: доказано, что всегда существует язык, который может быть задан LL(k)-грамматикой для некоторого $k > 0$, но не может быть задан LL($k-1$)-грамматикой. Таким образом, все LL(k)-грамматики для всех k представляют определенный интерес (другое дело, что распознаватели для них при больших значениях k будут слишком сложны). Интересно, что проблема эквивалентности для двух LL(k)-грамматик разрешима.

С другой стороны, для LR(k)-грамматик, составляющих класс LR-грамматик, доказано, что любой язык, заданный LR(k)-грамматикой с $k > 1$, может быть задан LR(1)-грамматикой. То есть LR(k)-грамматики с $k > 1$ интереса не представляют. Однако доказательство существования LR(1)-грамматики вовсе не означает, что такая грамматика всегда может быть построена (проблема преобразования КС-грамматик неразрешима).

На рис. 12.11 условно показана связь между некоторыми классами КС-грамматик, упомянутых в данном пособии. Из этой схемы видно, например, что любая LL-грамматика является LR-грамматикой, но не всякая LL-грамматика является LR(1)-грамматикой.

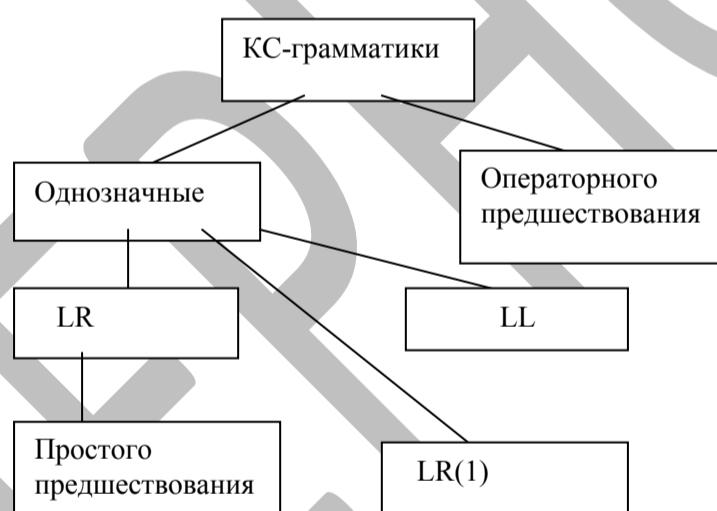


Рис. 12.11. Схема взаимосвязи некоторых классов КС-грамматик

Если вспомнить, что любой детерминированный КС-язык может быть задан, например, LR(1)-грамматикой, но в то же время, классы левоанализируемых и правоанализируемых грамматик несопоставимы, то напрашивается вывод: один и тот же детерминированный КС-язык может быть задан двумя или более несопоставимыми между собой грамматиками. Таким образом, можно вернуться к мысли о том, что проблема преобразования КС-грамматик неразрешима (на самом деле, конечно, наоборот: из неразрешимости проблемы преобразования КС-грамматик следует возможность задать один и тот же КС-язык двумя несопоставимыми грамматиками). Это, наверное, самый интересный вывод, который можно сделать из сопоставления разных классов КС-грамматик.

Отношения между классами КС-языков

КС-язык называется языком некоторого класса КС-языков, если он может быть задан КС-грамматикой из данного класса КС-грамматик. Например, класс LL-языков составляют все языки, которые могут быть заданы с помощью LL-грамматик.

Соотношение классов КС-языков представляет определенный интерес, оно не совпадает с соотношением классов КС-грамматик. Это связано с многократно уже упоминавшейся проблемой преобразования грамматик. Например, выше уже говорилось о том, что любой LL-язык является и LR(1)-языком — то есть язык, заданный LL-грамматикой, может быть

задан также и LR(1)-грамматикой. Однако не всякая LL-грамматика является при этом LR(1)-грамматикой и не всегда можно найти способ, как построить LR(1)-грамматику, задающую тот же самый язык, что и исходная LL-грамматика.

На рис. 12.12 приведено соотношение между некоторыми известными классами КС-языков [6, т. 2, 42, 47].

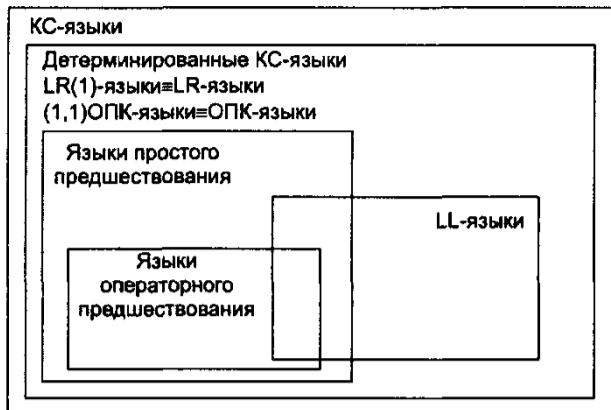


Рис. 12.12. Соотношение между различными классами КС-языков

Следует обратить внимание прежде всего на то, что интересующий разработчиков компиляторов в первую очередь класс детерминированных КС-языков полностью совпадает с классом LR-языков и, более того, совпадает с классом LR(1)-языков. То есть доказано, что для любого детерминированного КС-языка существует задающая его LR(1)-грамматика. Этот факт уже упоминался выше. Проблема состоит в том, что не всегда возможно найти такую грамматику и нет формализованного алгоритма, как ее построить в общем случае.

Также уже упоминалось, что LL-языки являются собственным подмножеством LR-языков: всякий LL-язык является одновременно LR-языком, но существуют LR-языки, которые не являются LL-языками. Поэтому LL-языки образуют более узкий класс, чем LR-языки.

Языки простого предшествования, в свою очередь, также являются собственным подмножеством LR-языков, а языки операторного предшествования — собственным подмножеством языков простого предшествования. Интересно, что языки операторного предшествования представляют собой более узкий класс, чем языки простого предшествования.

В то же время языки простого предшествования и LL-языки несопоставимы между собой: существуют языки простого предшествования, которые не являются LL-языками, и в то же время существуют LL-языки, которые не являются языками простого предшествования. Однако существуют языки, которые одновременно являются и языками простого предшествования, и LL-языками. Аналогичное замечание относится также к соотношению между собой языков операторного предшествования и LL-языков.

Можно еще отметить, что язык арифметических выражений над символами a и b , заданный грамматикой $G(\{+,-,/,*,a,b\}, \{S,T,E\}, P, S)$, $P = \{S \rightarrow S + T | S - T | T, T \rightarrow T * E | T / E | E, E \rightarrow (S) | a | b\}$, который многократно использовался в примерах в данном учебном пособии, подпадает под все указанные выше классы языков. Из приведенных ранее примеров можно заключить, что этот язык является и LL-языком, и языком операторного предшествования, а следовательно, и языком простого предшествования и, конечно, LR(1)-языком. В то же время этот язык по мере изложения материала пособия описывался различными грамматиками, не все из которых могут быть отнесены в указанные классы. Более того, он может быть задан с помощью грамматики, которая не являлась даже однозначной.

Таким образом, соотношение классов КС-языков не совпадает с соотношением задающих их классов КС-грамматик. Это связано с неразрешимостью проблем преобразования и эквивалентности грамматик, которые не имеют строго формализованного решения.

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 $A \rightarrow \epsilon$, 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 $P_1 = \{S \rightarrow 0A1, 0A \rightarrow 00A1, A \rightarrow \epsilon\}$ 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. Методы описания перевода

Грамматики описывают только синтаксические аспекты формальных языков. При построении трансляторов необходимо каждой входной цепочке (программе на входном языке) поставить в соответствие другую цепочку (программу на выходном языке). Отношение между входными и выходными цепочками называют переводом.

Возникает проблема (аналогичная описанию синтаксиса языка) - задания бесконечного перевода конечными средствами.

5.1 Перевод и семантика

Существуют два подхода к описанию перевода:

- использование системы, аналогичной формальной грамматике или являющейся её расширением, которая порождает перевод (точнее пару цепочек, принадлежавших переводу);
- применение автомата, распознающего входную цепочку и выдающего на выходе перевод этой цепочки.
- Естественными требованиями к устройству, выполнявшему трансляцию, являются:
- простота анализа и синтеза: мы должны, с одной стороны, по описанию перевода достаточно легко установить, какие пары цепочек принадлежат переводу, и, с другой стороны, построение описания перевода не должно быть непосильной задачей для разработчика;
- возможность автоматического построения транслятора по описанию перевода;
- небольшой объем и эффективность трансляции;
- корректность и т. п.

Перейдём к рассмотрению простейших способов описания перевода.

Определение. Пусть Σ — входной алфавит и Δ — выходной алфавит. Переводом с языка $L_1 \subseteq \Sigma^*$ на язык $L_2 \subseteq \Delta^*$ называется отношение τ из Σ^* в Δ^* , для которого L_1 — область определения, а L_2 — множество значений.

Если $(x, y) \in \tau$, то Цепочка y называется выходом для цепочки x . В общем случае в переводе τ для одной входной цепочки может быть задано более одной выходной цепочки.

Для языков программирования перевод должен быть функцией, т. е. для каждой входной программы должно быть не более одной выходной программы.

Простейшие переводы можно задать при помощи гомоморфизма.

Пример 5.11. Пусть требуется перевести символы, образующие целые числа со знаком, в имена. В этом случае гомоморфизм h можно определить следующим образом:

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, -, +\};$$

$$h(a) = a, \text{ если } a \in \Sigma;$$

Для $a \in \Sigma$ Гомоморфизм $h(a)$ определяется в соответствии с табл. 5.1.

Таблица 5.1. Определение гомоморфизма $h(a)$

a	h(a)	a	h(a)	a	h(a)
1	один	5	пять	9	девять
2	два	6	шесть	0	нуль
3	три	7	семь	+	плюс
4	четыре	8	восемь	-	минус

Используя данное определение, перевод входной цепочки "-15", содержащей запись числа -15, будет иметь вид: "минус один пять".

Гомоморфизм позволяет описывать только простейшие переводы. Даже такие простые переводы, как $\tau_1 = \{(x, y) \mid x - \text{ инфиксное выражение и } y - \text{ префиксная запись выражения } x\}$ и $\tau_2 = \{(x, y) \mid x - \text{ инфиксное выражение и } y - \text{ постфиксная запись выражения } x\}$ нельзя задать при помощи гомоморфизма, нужны более мощные формализмы.

5.2 СУ-схемы

Схема синтаксически управляемого (СУ-схема) перевода представляет собой систему, порождающую пары цепочек, принадлежавших переводу. Неформально схема синтаксически управляемого перевода представляет собой грамматику, в которой каждому правилу приписан элемент перевода. При порождении цепочки каждый раз, когда правило участвует в этом процессе, элемент перевода генерирует часть выходной цепочки, соответствующей части входной цепочки, порождённой этим правилом.

Рассмотрим схему синтаксически управляемого перевода, описывающую перевод скобочного

арифметического выражения, порождаемого грамматикой G0, в соответствующую польскую инверсную запись. Схема перевода представлена в табл. 5.2.

Таблица 11.2. СУ-схема перевода

	Правило грамматики	Элемент перевода
(1)	$E \rightarrow E + T$	$E = ET +$
(2)	$E \rightarrow T$	$E = T$
(3)	$T \rightarrow T * P$	$T = TP *$
(4)	$T \rightarrow P$	$T = P$
(5)	$P \rightarrow (E)$	$P = E$
(6)	$P \rightarrow i$	$P = i$

В этой схеме правилу Грамматики $E \rightarrow E + T$ соответствует элемент перевода $E = ET+$, который интерпретируется следующим образом: "Перевод нетерминала E , стоявшего в левой части элемента перевода, представляет собой перевод, порождённый символом E , стоявшим в правой части этого элемента, за которым следуют переводы, порождаемые символом T , и символ '+'".

Используя СУ-схему, приведённую в табл. 5.2, определим перевод цепочки $a^*(b + c)$.
Замечание: для указания местоположения символов выходной цепочки, являющихся переводом символов входной цепочки, будем считать символы a, b, c значениями терминального символа i грамматики.

В начале определим левый вывод входной Цепочки:

$$\begin{aligned} E &\Rightarrow_{(2)} T \Rightarrow_{(3)} T * P \Rightarrow_{(4)} P * P \Rightarrow_{(6)} a^* P \Rightarrow_{(5)} a^* (E) \Rightarrow_{(1)} a^* (E + T) \Rightarrow_{(2)} \\ &a^* (T + T) \Rightarrow_{(4)} a^* (P + T) \Rightarrow_{(6)} a^* (b + T) \Rightarrow_{(4)} a^* (b + P) \Rightarrow_{(6)} \\ &a^* (b + c). \end{aligned}$$

Полученный вывод 23465124646 определяет последовательность использовав- неё элементов перевода при порождении выходной Цепочки, поэтому последовательность выводимых пар цепочек имеет следующий вид:

$$\begin{aligned} (E, E) &\Rightarrow_{(2)} (T, T) \\ &\Rightarrow_{(3)} (T * P, TP *) \\ &\Rightarrow_{(4)} (P * P, PP *) \\ &\Rightarrow_{(6)} (a^* P, aP *) \\ &\Rightarrow_{(5)} (a^* (E), a^* E) \\ &\Rightarrow_{(1)} (a^* (E + T), a^* ET +) \\ &\Rightarrow_{(2)} (a^* (T + T), a^* TT +) \\ &\Rightarrow_{(4)} (a^* (P + T), a^* PT +) \\ &\Rightarrow_{(6)} (a^* (b + T), a^* bT +) \\ &\Rightarrow_{(4)} (a^* (b + P), a^* bP +) \\ &\Rightarrow_{(6)} (a^* (b + c), a^* bc +). \end{aligned}$$

Каждая выходная цепочка при порождении части входной цепочки получается путем замены (подходившего) не терминала выходной цепочки путем замены (подходившего) не терминала выходной Цепочки значением соответствующего ему элемента перевода.

Рассмотренная схема перевода относится к классу схем, называемых схемами синтаксически управляемого перевода.

Определение: схемой синтаксически управляемого перевода называется пятёрка $T = (N, \Sigma, \Delta, R, S)$, где:

- N — конечное множество нетерминальных символов;
- Σ — конечный входной алфавит;
- Δ — конечный выходной алфавит;
- R — конечное множество правил вида $A \rightarrow \alpha, \beta$, где $\alpha \in (N \cup \Sigma)^*$, $\beta \in (N \cup \Delta)^*$ и вхождения не терминалов в цепочку β образуют перестановку вхождений не терминалов в цепочку α ;
- S — начальный символ ($S \in N$).

По определению, каждому вхождению определённого не терминала в цепочку α соответствует некоторое вхождение этого не терминала в цепочку β . Если некоторый не терминал входит в цепочку α более одного раза, то может возникнуть неоднозначность. Для исключения этого будем пользоваться верхними целочисленными индексами, например в правиле $A \rightarrow B^{(1)}CB^2, B^2B^{(1)}C$ первой, второй и третьей позиции цепочки $B(1)CB(2)$ соответствуют вторая, третья и первая позиции цепочки $B^2B^{(1)}C$.

Примечание: выводимая пара цепочек СУ-схемы T определяется рекурсивно следующим образом:

1) (S, S) — выводимая пара, в которой символы S соответствуют друг другу.

2) Если $(\alpha A \beta, \alpha' A \beta')$ — выводимая пара, в которой два выделенных вхождения

не терминала A соответствуют друг другу, и $A \rightarrow \alpha, \beta$ — правило из R , то $(\alpha\gamma\beta, \alpha'\gamma'\beta')$ — выводимая пара.

Вхождения не терминалов в цепочки γ и γ' соответствуют Друг Другу точно так же, как они соответствовали в правиле СУ-схемы. Вхождения не терминалов в цепочки α и β соответствуют вхождениям не терминалов в цепочки α' и β' в новой выводимой паре точно так же, как они соответствовали в старой выводимой паре. При необходимости это соответствие будет указываться верхними индексами.

Если между парами цепочек $(\alpha A \beta, \alpha' A \beta')$ и $(\alpha \gamma \beta, \alpha' \gamma' \beta')$ установлена описанная ранее связь, то будем писать $(\alpha A \beta, \alpha' \gamma' \beta') \Rightarrow_T (\alpha \gamma \beta, \alpha' \gamma' \beta')$.

Транзитивное замыкание, рефлексивно-транзитивное замыкание и k -Ю степень отношения \Rightarrow_T будем обозначать как \Rightarrow_T^+ , \Rightarrow_T^* , и \Rightarrow_T^k , соответственно. В очевидных случаях можно опускать индекс 'T'.

Переводом, определяемым схемой T , или $\square(T)$ называют множество пар цепочек $\{(x, y) | (S, S) \Rightarrow^*(x, y), x \in \Sigma^* \text{ и } y \in \Delta^*\}$.

Пример 5.2

Пусть СУ-схема задана следующим образом f31:

$$T = (\{S\}, \{i, +\}, \{i, +\}, R, S),$$

где множество правил R содержит правила

$$(1) \quad S \rightarrow + S(1) S(2), S(1) + S(2)$$

$$(2) \quad S \rightarrow i, i$$

Рассмотрим вывод в данной СУ-схеме:

$$\begin{aligned} (S, S) &\Rightarrow_{(1)} (+S^{(1)}S^{(2)}) \\ &\Rightarrow_{(1)} (+ + S^{(3)} S^{(4)} S^{(2)}) \\ &\Rightarrow_{(1)} (+ + + S^{(5)} S^{(6)} S^{(4)} S^{(2)}) \\ &\Rightarrow_{(2)} (+ + + a S^{(6)} S^{(4)} S^{(2)}) \\ &\Rightarrow_{(2)} (+ + + a b S^{(4)} S^{(2)}) \\ &\Rightarrow_{(2)} (+ + + a b c S^{(2)}) \\ &\Rightarrow_{(2)} (+ + + a b c d) \end{aligned} \quad , \quad \begin{aligned} S^{(1)+} S^{(2)} \\ S^{(3)} + S^{(4)} + S^{(2)} \\ S^{(5)} + S^{(6)} + S^{(4)} + S^{(2)} \\ a + S^{(6)} S^{(4)} + S^2 \\ a + b + S^{(4)} + S^{(2)} \\ a + b + c + S^{(2)} \\ a + b + c + d. \end{aligned}$$

Входной Цепочке $+++$ соответствует выходная Цепочка $i+i+i+i$.

Очевидно, что перевод, определяемый СУ-схемой, имеет вид:

$$\tau(T) = \{(x, i(+i)k) | k \geq 0 \text{ и } x - \text{ префиксная запись выражения } i(+i)k\}.$$

Примечание: если $T = (T, \Sigma, \Delta, R, S)$ — СУ-схема, то $\tau(T)$ называется синтаксически управляемым переводом.

Примечание: $G_i = (N, \Sigma, P, S)$, где $P = \{A \rightarrow \alpha \mid A \rightarrow \alpha, \beta, \in R\}$, называется входной грамматикой СУ-схемы t .

Примечание: $G_0 = (N, \Delta, P', S)$, где $P' = \{A \rightarrow \beta \mid A \rightarrow \alpha, \beta, \in R\}$, называется выходной грамматикой СУ-схемы t .

К достоинствам СУ-схемы относится не только ее простота, но и нагларь- носить ее интерпретации.

Синтаксически управляемый перевод можно трактовать как метод преобразования деревьев выводов входной грамматики G_i в деревья выводов выходной грамматики G_0 . Перевод данной входной цепочки и x можно получить, построив ее дерево вывода, затем преобразовав это дерево в дерево вывода в выходной грамматике и, наконец, взяв крону выходного дерева в качестве перевода цепочки x [3].

Алгоритм 5.1 Преобразование деревьев при помощи СУ-схемы

Вход: СУ-схема $T = (N, \Sigma, \Delta, R, S)$ с входной грамматикой $G_i = (N, \Sigma, P, S)$ и выходной грамматикой $G_0 = (N, \Delta, P', S)$ и дерево вывода D в G_i , с кроной, принадлежавшей Σ^* .

Выход: Некоторое дерево вывода D' в G_0 , такое, что если x и y — кроны деревьев D и D' соответственно, то $(x, y) \in \tau(T)$

Описание алгоритма:

1. Пусть Данный шаг применяется к внутренней вершине n Дерева D , имею шей к прямых потомков n_1, \dots, n_k .

1.1 УстраниТЬ из множества вершин n_1, \dots, n_k листья, помеченные терминальными символами или ϵ .

1.2 Пусть $A \rightarrow \alpha$ — правило входной грамматики G_i , соответствую шее вершине n и ее прямым потомкам, т. е. A — метка вершины n , и α образуется конкатенацией меток вершин n_1, \dots, n_k . Выбрать из R некоторое правило вида $A \rightarrow \alpha, \beta$ (если таких правил несколько, то выбор произволен).

Переставить оставшиеся прямые потомки вершины n (если они есть) согласно соответству между вхождениями не терминалов в α и β . (Поддеревья, корнями которых служат эти потомки, переставляются вместе с ними).

1.3 Добавить в качестве прямых потомков вершины n листья с метками так, чтобы метки всех ее прямых потомков образовали цепочку β .

1.4 Применить Шаг 1 к прямым потомкам вершины n , не являвшимся листьями, в порядке слева направо.

2. Повторять Шаг 1 рекурсивно, начиная с корня дерева D . Результатирующим деревом будет D' .

Пример 5.3

Рассмотрим СУ-схему $1 = (\{S, A\}, \{0, 1\}, \{a, b\}, R, S)$, где R состоит из правил:

- (1) $S \rightarrow 0AS, SAa$
- (2) $A \rightarrow 0SA, ASA$
- (3) $S \rightarrow 1, b$
- (4) $A \rightarrow 1, b$

Во входной грамматике этой схемы цепочка 0010111 имеет вывод 1232343 (дерево вывода приведено на

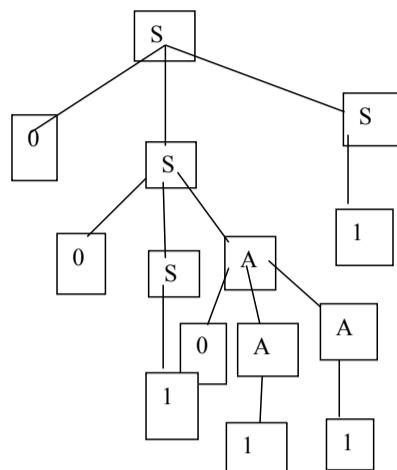


рис. 5.1).

Рис. 5.1. Дерево вывода цепочки 0010111

Преобразуем это дерево при помощи алгоритма 5.1.

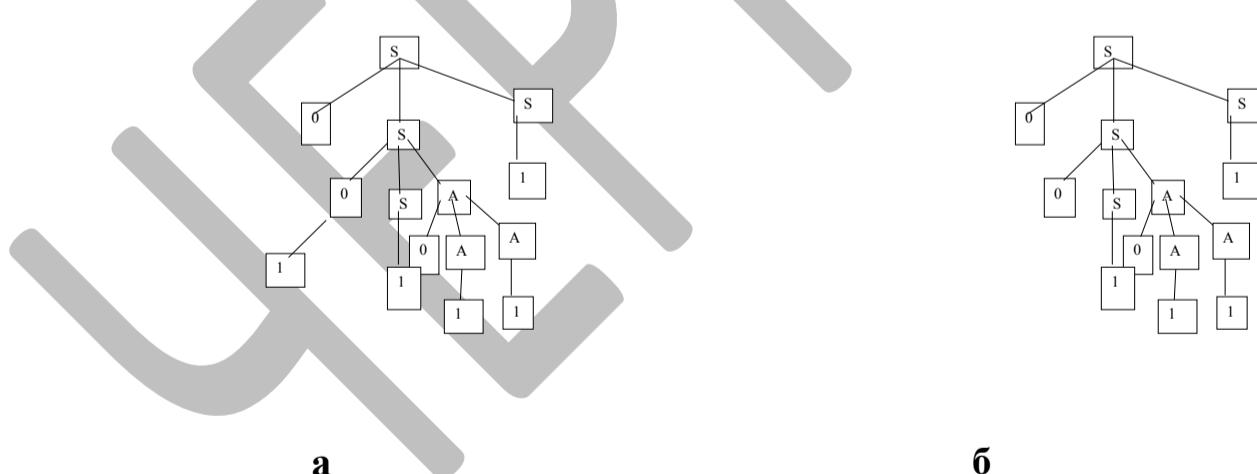
Шаг 1. Применим шаг 1 алгоритма к корню Дерева S:

Устранием левый лист дерева, помеченный 0. Находим правило входной грамматики (1) $S \rightarrow 0AS, SAa$, соответствующее корню дерева S. У этого правила только один элемент перевода SA, который определяет, что нужно поменять местами прямые потомки корня A и S. Добавляем прямой потомок к корню дерева и помечаем его a.

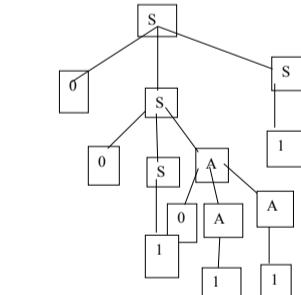
результат этих действий приведена на рис. 5.2, а.

Шаг 1. Рассматриваем прямой потомок корня дерева – вершину, помеченную символом A. Устранием листа 0. Используя правило (2) СУ-схемы, меняем местами прямые потомки рассматриваемой вершины. Добавляем прямой потомка к рассматриваемой вершине и получаем дерево, приведённое на рис. 5.2, б.

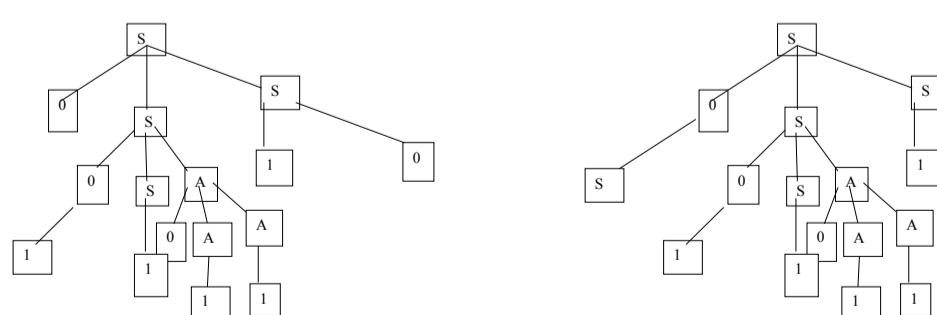
Шаг 1. Рассматриваем прямой потомок корня дерева – вершину, помеченную символом S. Устранием листа 1. Находим соответствующее правило СУ-схемы – правило номер (3), и так как потомков у рассматриваемой вершины нет, то создаём прямой потомок b. Для рассматриваемой вершины и получаем Дерево, приведённое на рис. 5.2, в.



а



б



в

г

Рис. 5.2. Преобразование деревьев

Продолжая рекурсивно выполнение шага 1 алгоритма, получаем результирующее дерево, приведённое на рис. 5.2, г.

В [3] Доказаны следующие утверждения:

- 1) Если x и y – кроны деревьев D и D' из алгоритма 5.1, то $(x, y) \in \tau(T)$.
- 2) Если $(x, y) \in \tau(T)$, то существуют Дерево вывода D с кроной x и такая последовательность выборов вершин при обращении к Шагу 2.2 алгоритма, что в результате получается Дерево D' с кроной y .

Замечание: порядок применений шага 2 алгоритма 5.1 к вершинам дерева не важен. Можно выбрать любой порядок, при котором каждая внутренняя вершина рассматривается точно один раз.

Примечание: СУ-схема $T = (N, \Sigma, \Delta, R, S)$ называется простой, если для каждого правила $A \rightarrow \alpha, \beta \in R$ соответствующие друг другу вхождения не терминалов встречаются в α и β в одном и том же порядке.

Примечание: перевод, определяемый простой СУ-схемой, называется простым синтаксически управляемым переводом.

Соответствие нетерминалов в выводимой паре простой СУ-схемы определяется порядком, в котором эти нетерминалы появляются в цепочках, поэтому для нее легко построить транслятор, представляющий собой преобразователь с магазинной памятью.

По этой причине простые СУ-переводы образуют важный класс переводов.

Пример 5.4

Рассмотрим простую СУ-схему, имеющую следующие правила:

- (1) $E \rightarrow (E), E$
- (2) $E \rightarrow E + E, E + E$
- (3) $E \rightarrow T, T$
- (4) $T \rightarrow (T), T$
- (5) $T \rightarrow A * A, A * A$
- (6) $T \rightarrow i, i$
- (7) $A \rightarrow (E + E), (E + E)$
- (8) $A \rightarrow T, T$

и вывод входной цепочки $((i * (i * i) + i) * i)$, равный 3457358684586863686.

Соответствующий вывод пар цепочек имеет вид:

$$\begin{array}{ll}
 (E, E) \Rightarrow_{(3)} (T, T) \\
 \Rightarrow_{(4)} ((T, T)) \\
 \Rightarrow_{(5)} ((A * A), A * A) \\
 \Rightarrow_{(7)} (((E + E) * A), (E + E) * A) \\
 \Rightarrow_{(3)} (((T + E) * A), (T + E) * A) \\
 \Rightarrow_{(5)} (((A * A + E) * A), (A * A + E) * A) \\
 \Rightarrow_{(8)} (((T * A + E) * A), (T * A + E) * A) \\
 \Rightarrow_{(6)} (((i * A + E) * A), (i * A + E) * A) \\
 \Rightarrow_{(8)} (((i * T + E) * A), (i * T + E) * A) \\
 \Rightarrow_{(4)} (((i * (T) + E) * A), (i * T + E) * A) \\
 \Rightarrow_{(5)} (((i * (A * A) + E) * A), (i * A * A + E) * A) \\
 \Rightarrow_{(8)} (((i * (T * A) + E) * A), (i * T * A + E) * A) \\
 \Rightarrow_{(6)} (((i * (i * A) + E) * A), (i * i * A + E) * A) \\
 \Rightarrow_{(8)} (((i * (i * T) + E) * A), (i * i * T + E) * A) \\
 \Rightarrow_{(6)} (((i * (i * i) + E) * A), (i * i * i + E) * A) \\
 \Rightarrow_{(3)} (((i * (i * i) + T) * A), (i * i * i + T) * A) \\
 \Rightarrow_{(6)} (((i * (i * i) + i) * A), (i * i * i + i) * A) \\
 \Rightarrow_{(8)} (((i * (i * i) + i) * T), (i * i * i + i) * T) \\
 \Rightarrow_{(6)} (((i * (i * i) + i) * i), (i * i * i + i) * i).
 \end{array}$$

Для входной цепочки $((i * (i * i) + i) * i)$ СУ-схема порождает выходную цепочку $(i * i * i + i) * i$. анализ позволяет сделать вывод, что СУ-схема отображает арифметические выражения, порождаемые грамматикой G_0 , в арифметические выражения, не содержащие избыточных скобок.

Рассмотрение СУ-схем позволяет понять идеи, лежащие в основе методов описания синтаксически управляемого перевода: включение в грамматику, описывающую порождение цепочек языка, элементов перевода, позволяющие описывать порождение выходных Цепочек.

Для получения дополнительных сведений о СУ-схемах можно обратиться к [3, 4].

5.3. Транслирующие грамматики

Рассматривая процесс перевода инфиксных арифметических выражений в польскую инверсионную запись, попытаемся построить некоторый процессор, выполняющий этот перевод.

Если на входной ленте процессора находится цепочка $a + b * c$, то процессор должен работать следующим образом:

- прочитать входной символ, a ;
- выдать символ, a на выходную ленту;
- прочитать входной символ $'+'$;
- прочитать входной символ Б ;
- выдать символ Б на выходную ленту;
- прочитать входной символ $'*'$;
- прочитать входной символ c ;

- выдать символ с на выходную ленту;
- выдать символ '*' на выходную ленту;
- выдать символ '+' на выходную ленту.

Этот план работы определяется тем, что порядок следования операндов в инфиксной записи и ПОЛИЗ совпадает, а операции в ПОЛИЗ должны следовать сразу за своими operandами. Если операцию Чтения символа из входной ленты мы будем обозначать Читаемыми символами, а операцию записи — символами, помещаемыми на выходную ленту, заключёнными в фигурные скобки (операционными символами), то работа по переводу может быть записана следующим образом:

а {a} + b {b} * c {c} {*} {+}.

Транслирующая грамматика – это КС-грамматика, множество терминалов которой разбито на два множества: множество входных и множество операционных символов.

Дадим формальное определение транслирующей грамматики.

Определение: транслирующей грамматикой называется пятерка объектов $G^T = (N, \Sigma_i, \Sigma_a, P, S)$, где — словарь входных символов, Σ_a — словарь операционных символов, N — нетерминальный словарь, $S \in N$ — начальный символ транслирующей грамматики, P — конечное множество правил вывода вида $A \rightarrow \alpha$, в которых $A \in N$, а $\alpha \in (\Sigma_i \cup \Sigma_a \cup N)^*$.

Транслирующая грамматика G_0^T , описывающая перевод инфиксных арифметических выражений в ПОЛИЗ, содержит следующие правила:

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

Грамматика, полученная из транслирующей грамматики путем вычеркивания всех операционных символов, называется входной грамматикой для этой транслирующей грамматики. Язык, порождаемый входной грамматикой, называется входным языком.

Например, входной грамматикой G_0 для транслирующей грамматики G_0^T является КС-грамматика для инфиксных арифметических выражений, содержащая правила:

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

Рассмотрим левый вывод цепочки $i + i * i$ во входной грамматике G_0 :

$$E \Rightarrow E + T \Rightarrow T + T \Rightarrow P + T \Rightarrow i + T \Rightarrow i + T * P \Rightarrow i + P * P \Rightarrow i + i * P \Rightarrow i + i * i.$$

Применив эту же последовательность правил к соответствующим нетерминалам транслирующей грамматики G_0^T получим следующий левый вывод:

$$E \Rightarrow E + T \{+\} \Rightarrow T + T \{+\} \Rightarrow P + T \{+\} \Rightarrow i \{i\} + T \{+\} \Rightarrow i \{i\} + T * P \{*\} \{+\} \Rightarrow i \{i\} + P * P \{*\} \{+\} \Rightarrow i \{i\} + i \{i\} * P \{*\} \{+\} \Rightarrow i \{i\} + i \{i\} * i \{i\} \{*\} \{+\}.$$

Цепочки языка, порождаемого транслирующей Грамматикой, называют активными цепочками. Входной частью активной цепочки называется последовательностью входных символов, полученная из активной цепочки после удаления всех операционных символов. Операционной частью активной цепочки называется последовательность операционных символов, полученная путем вычёркивания из активной цепочки всех входных символов.

Например, последовательность $i + i * i$ является входной Частью активной цепочки $i \{i\} + i \{i\} * i \{i\} \{*\} \{+\}$, а последовательностью $\{i\} \{i\} \{i\} \{*\} \{+\}$ – ее операционной частью.

Примечание: множество всех пар, первыми элементами которых являются входные части активной цепочки, а вторыми элементами — операционные части активной цепочки, называется синтаксически управляемым переводом $\tau(G^T)$, определяемым транслирующей грамматикой G^T .

В восходящих методах обработки языков широко применяются постфиксные транслирующие грамматики.

Примечание: транслирующая грамматика называется постфиксной транслирующей грамматикой тогда и только тогда, когда все операционные символы в правых частях правил вывода расположены правее всех входных и нетерминальных символов.

Примером постфиксной транслирующей Грамматики является Грамматика G_0^T .

Для получения транслирующей грамматики входной язык описывается входной КС-грамматикой. Затем в правила вывода этой грамматики вставляются операционные символы для описания семантических действий, связанных с соответствующими правилами.

Рассмотрим получение транслирующей грамматики G_0^T по входной Грамматике G_0 .

В ПОЛИЗ operandы располагаются в той же последовательности, что и в инфиксной записи, а знаки операций следуют непосредственно после своих operandов в том порядке, в котором они выполняются. Для того чтобы идентификатор i выдавался сразу после его прочтения, правило (6) Грамматики G_0 преобразуется к виду $P \rightarrow i \{i\}$. Чтобы знак операции сложения выдавался непосредственно после своих operandов, правило (1) заменяется на $E \rightarrow E + T \{+\}$. Это правило интерпретируется следующим образом: обработка

арифметического выражения Е состоит из обработки первого операнда Е, чтения символа '+', обработки второго операнда Т и выдачи символа '+'. Аналогичные рассуждения позволяют следующим образом преобразовать правило (3) грамматики:

$T \rightarrow T * P \{*\}$. Один и тот же перевод можно определить разными транслирующими грамматиками

Пример 5.5

Пусть две транслирующие грамматики G_1^T и G_2^T , определяющие перевод в ПОЛИЗ инфиксных бесскобочных арифметических выражений, выполняемых в порядке написания операций, выглядят следующим образом:

G_1^T

- (1) $E \rightarrow E + T \{+\}$
- (2) $E \rightarrow E * T \{*\}$
- (3) $E \rightarrow T$
- (4) $T \rightarrow i \{i\}$

G_2^T

- (1) $E \rightarrow i \{i\} R$
- (2) $R \rightarrow + i \{i\} \{+\} R$
- (3) $R \rightarrow * i \{i\} \{*\} R$
- (4) $R \rightarrow \epsilon$

Рассмотрим выполнение перевода цепочки $i + i * i$ с использованием транслирующих грамматик G_1^T и G_2^T .

При использовании Грамматики G_1^T активная цепочка порождается следующим образом:

$$\begin{aligned} E &\Rightarrow_{(2)} (2) E * T \{*\} \\ &\Rightarrow_{(1)} E + T \{+\} * T \{*\} \\ &\Rightarrow_{(3)} T + T \{+\} * T \{*\} \\ &\Rightarrow_{(4)} i \{i\} + T \{+\} * T \{*\} \\ &\Rightarrow_{(4)} i \{i\} + i \{i\} \{+\} * T \{*\} \\ &\Rightarrow_{(4)} i \{i\} + i \{i\} \{+\} * i \{i\} \{*\} \end{aligned}$$

Полученная операционная часть цепочки $\{i\} \{i\} \{+\} \{i\} \{*\}$ является переводом входной цепочки $i + i * i$.

Использование второй Грамматики G_2^T позволяет получить следующее порождение активной Цепочки:

$$\begin{aligned} E &\Rightarrow_{(1)} i \{i\} R \\ &\Rightarrow_{(2)} i \{i\} + i \{i\} \{+\} R \\ &\Rightarrow_{(3)} i \{i\} + i \{i\} \{+\} * i \{i\} \{*\} R \\ &\Rightarrow_{(4)} i \{i\} + i \{i\} \{+\} * i \{i\} \{*\} \end{aligned}$$

Полученная операционная часть цепочки $\{i\} \{i\} \{+\} \{i\} \{*\}$ также является переводом входной цепочки. На рис. 5.3 приведены деревья вывода Цепочки $i \{i\} + i \{i\} \{+\} * i \{i\} \{*\}$ с использованием грамматик G_1^T и G_2^T соответственно.

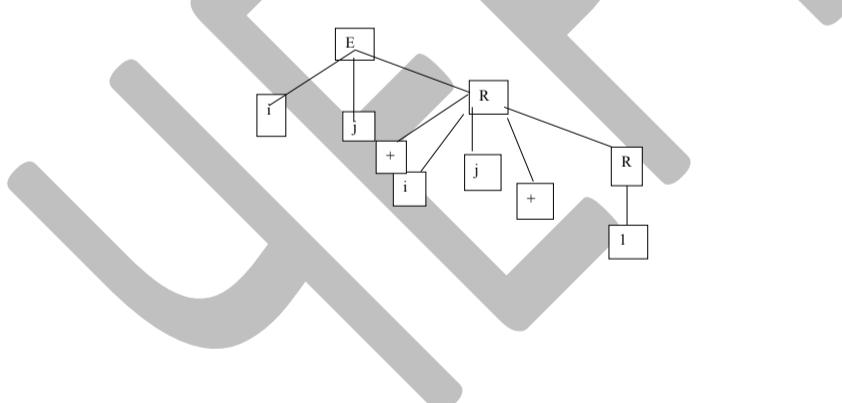


Рис. 5.3. Деревья вывода цепочки $i \{i\} + i \{i\} \{+\} * i \{i\} \{*\}$ в грамматиках G_1^T

Во многих практических приложениях появление входного символа в активной цепочке можно интерпретировать как обозначение операции чтения этого символа некоторым устройством, а вхождению операционного символа в правило вывода можно сопоставить операцию выдачи символа, заключённого в фигурные скобки. При таком подходе входную часть активной цепочки называют входной цепочкой, операционную часть – входной Цепочкой, а транслирующую грамматику, определяющую перевод, грамматикой цепочечного перевода.

Примером грамматики цепочечного перевода является грамматика G_0^T :

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

Другая возможная интерпретация операционных символов - имена семантических процедур, вызываемых при обработке входной цепочки. В этом случае активная цепочка определяет последовательность вызовов семантических процедур и моменты времени, в которые выполняются эти вызовы по отношению к моментам

чтения входных символов.

5.4 Атрибутные транслирующие грамматики

При определении транслирующей грамматики в понятие входного символа не включалось представление о том, что он является лексемой и состоит из двух компонентов: типа и значения. На самом деле транслирующая грамматика способна описывать перевод только той части символа, которая задаёт его тип. Рассмотрим, как можно расширить понятие грамматики цепочечного перевода, чтобы использовать при переводе оба компонента входного символа. Расширенная транслирующая грамматика получила название атрибутной транслирующей грамматики (АТ-Грамматики). АТ-Грамматики были предложены Дональдом Кнутом и в дальнейшем изучались рядом учёных. Материал данного раздела в основном базируется на результатах, приведённых в [2, 27, 28].

В АТ-Грамматике символы снабжаются атрибутами, которые могут принимать значения из некоторого множества допустимых значений и интересуются как семантическая информация, связанная с конкретным вхождением символа в правило вывода грамматики. При таком подходе значения атрибутов связываются с вершинами дерева вывода в АТ-грамматике, а правила вычисления значений атрибутов сопоставляются правилам вывода грамматики.

Все атрибуты нетерминальных и операционных символов делятся на синтезированные и унаследованные атрибуты.

5.4.1 Синтезированные атрибуты

Рассмотрим получение АТ-Грамматики из исходной транслирующей грамматики.

Пример 5.6

Пусть задана транслирующая грамматика G_3^T , допускающая в качестве входа арифметические выражения, построенные из символов множества

$\{c, +, *, (,)\}$, где c - лексема, является целочисленной константой и порождающая на выходе значение этого выражения:

- | | |
|---------------------------------|-------------------------|
| (O) $S \rightarrow E \{OTBET\}$ | |
| (1) $E \rightarrow E + T$ | (4) $T \rightarrow P$ |
| (2) $E \rightarrow T$ | (5) $P \rightarrow c$ |
| (3) $T \rightarrow T * P$ | (6) $P \rightarrow (E)$ |

На рис. 5.4, а изображено дерево вывода для входной цепочки $c_5 + c_2 * c_3$ (индексы обозначают значения лексем). Этой входной цепочке должна соответствовать выходная цепочка $\{OTBET\}$.

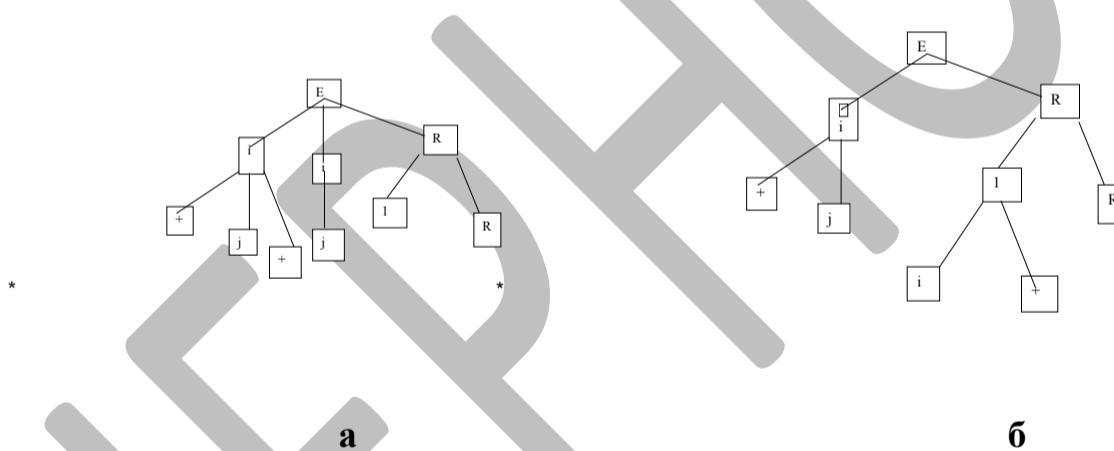


Рис. 5.4. Дерево вывода цепочки $c_5 + c_2 * c_3$ во входной (а) и атрибутной транслирующей (б) грамматиках.

Поскольку любое вхождение не терминалов E , T и P в дерево вывода представляет собой некоторое подвыражение входного выражения, введём для каждого из символов E , T и P по одному атрибуту, значением которого будет значение подвыражения, порождаемого этим не терминалом.

Входной символ c и операционный символ $\{OTBET\}$ также имеют по одному атрибуту. Значение атрибута символа c равно значению константы, которую он представляет, а значением символа $\{OTBET\}$ является результат выражения.

Все атрибуты символов могут принимать любые целочисленные значения из области допустимых для выражений, описываемых соответствующей входной грамматикой G_3 . На рис. 5.4, б изображено дерево вывода той же цепочки, что и на рис. 5.4, а, но на этом дереве не терминалы и операционный символ $\{OTBET\}$ помечены значениями своих атрибутов.

Выберем для каждого атрибута каждого вхождения символа в правило вывода грамматики уникальное имя и включим атрибуты в правила вывода в виде подстрочных индексов соответствующих символов.

Замечание: одни и те же имена атрибутов можно использовать в нескольких правилах, т. к. они локальны в правиле грамматики.

Сопоставим каждому правилу вывода Грамматики правило вычисления значения атрибута не терминала из левой Части правила, которому соответствует нетерминальная вершина дерева вывода.

Рассмотрим, например, вершину T дерева, изображённого на рис. 5.4, а. Правило вывода, применённое к этой вершине, имеет вид: $T \rightarrow T * P$. Оно определяет, что значение подвыражения, порождённого не терминалом из левой Части правила, равно значению подвыражений, порождённых символами P и T из правой Части правила, которые являются прямыми потомками вершины T . Если р

- атрибут символа Т из левой Части правила, а q и r - атрибуты символов Т и Р из правой части правила, то правило вычисления значения атрибута р будет иметь вид: $p \leftarrow q * r$, где символ ' \leftarrow ' - знак операции присваивания.

Таким образом, начиная с атрибутов входных символов и поднимаясь по Дереву от кроны к корню, можно определить все значения атрибутов Нетер-мина лов из левых Частей правил вывода. Правила вычисления значений атрибутов Для АТ-Грамматики G_3^T , полученной по транслирующей грамматике G_3^T , будут выглядеть следующим образом:

$$\begin{array}{ll}
 (O) \quad S \rightarrow E_q \{ОТВЕТ\} p & (4) \quad T_p \rightarrow * P_q \\
 & q \leftarrow p \\
 (1) \quad E_p \rightarrow E_q + T_r & (5) \quad P_p \rightarrow (E_q) \\
 & p \leftarrow q + r \\
 (2) \quad E_p \rightarrow & (6) \quad P_p \rightarrow C_q \\
 & p \leftarrow q \\
 (3) \quad T_p \rightarrow T_q * P_r & \\
 & p \leftarrow q * r
 \end{array}$$

Атрибуты, значения которых вычисляются при движении по дереву вывода снизу вверх, традиционно называют синтезированными. Термин "синтезированный" подчёркивает, что значение атрибута синтезируется из значений атрибутов потомков.

Рассмотренные ранее правила вычисления значений атрибутов не распространяются на атрибут операционного символа {ОТВЕТ}, который относится к классу унаследованных атрибутов.

5.4.2 Унаследованные атрибуты

Термин "унаследованный" означает то, что значение атрибута зависит от значений атрибутов предка символа или атрибутов его соседей в дереве вывода. Например, в грамматике G3T, значение атрибута г символа {ОТВЕТ} равно значению атрибута не терминала Е (соседа слева), порождающего все выражение. Оно может быть вычислено только после того, как будут определены значения всех синтезированных атрибутов не терминалов.

Рассмотрим ещё один пример АТ-Грамматики.

Пример 5.7

Пусть входная КС-грамматика G, порождающая описания переменных в некотором языке программирования, выглядит следующим образом:

- 1) $D \rightarrow t \ i \ L$
- 2) $L \rightarrow , \ i \ L$
- 3) $L \rightarrow \epsilon$

Множество входных символов этой грамматики состоит из лексем: запятая, идентификатор *i* и имя типа *t*.

Семантика обработки описания заключается в занесении типа переменной в определённое поле элемента таблицы идентификаторов. Эту операцию можно выполнить с помощью семантической процедуры УСТАНОВИТЬ_ТИП с двумя параметрами: указатель на элемент таблицы идентификаторов, соответствующий описываемой переменной, и тип переменной. Процедуру УСТАНОВИТЬ_ТИП лучше всего вызывать сразу после распознавания идентификатора. Указанная последовательность действий может быть описана транслирующей грамматикой G_T :

- 1) $D \rightarrow t \ i \ \{\text{ТИП}\} \ L$
- 2) $L \rightarrow , \ i \ \{\text{ТИП}\} \ L$
- 3) $L \rightarrow \epsilon$

В этой грамматике вызову процедуры УСТАНОВИТЬ_ТИП соответствует операционный символ {ТИП}.

Введём в транслирующую грамматику G^T атрибуты и правила их вычисления. Входные символы *t* и *i* имеют по одному атрибуту. Значением атрибута символа *i* является указатель на элемент таблицы идентификаторов, а атрибут символа *t* может принимать значения из множества {ЦЕЛЫЙ, ВЕЩЕСТВЕННЫЙ, ЛОГИЧЕСКИЙ}.

Операционный символ {ТИП} должен иметь два унаследованных атрибута, значения которых совпадают со значениями соответствующих фактических параметров процедуры УСТАНОВИТЬ_ТИП. Значения унаследованных атрибутов символа {ТИП} Для первого правила вывода приравниваются значениям соответствующих атрибутов входных символов *i* и *t*, входящих в правую часть того же правила левее символа {ТИП}.

Во второе правило тип описываемых переменных можно передать через унаследованный атрибут не терминала *L*. Значение этого унаследованного атрибута будет передаваться по дереву сверху вниз, начиная с вершины, где он получает начальное значение, равное значению атрибута входного символа *t*.

Атрибутная транслирующая Грамматика G^A выглядит следующим образом:

- (1) $D \rightarrow t_r \ i_a \ \{\text{ТИП}\}_{a1, r1} \ L_{r2}$

$$\begin{array}{l}
 a_1 \leftarrow a \\
 r_1, r_2 \leftarrow r \\
 (2) L_r \rightarrow , i_a \{\text{ТИП}\}_{a_1, r_1} L_{r_1} \\
 \quad \quad \quad a_1 \leftarrow a \\
 \quad \quad \quad r_1, r_2 \leftarrow r \\
 (3) L_r \rightarrow \varepsilon
 \end{array}$$

Замечание: запись атрибутного правила в виде $r_1, r_2 \leftarrow r$ означает, что значение r присваивается одновременно атрибутам r_1 и r_2 .

Входной цепочке $t_{\text{ЦЕЛЫЙ}} i_5, i_9, i_2$ соответствует дерево вывода грамматике G^A , изображённое на рис. 5.5.

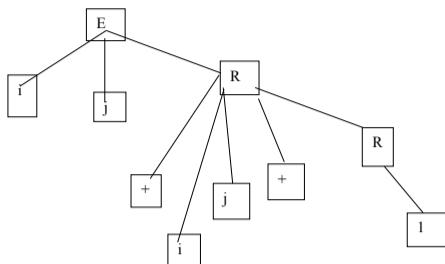


Рис. 5.5. Дерево вывода

В рассмотренных примерах все операционные символы имели унаследованные атрибуты. При определении атрибутной транслирующей грамматики можно обойтись без синтезированных атрибутов операционных символов. Необходимость в синтезированных атрибутах операционных символов может возникнуть в некоторых практических реализациях переводов, поэтому в определение АТ-Грамматики включены оба типа атрибутов.

5.4.3. Определение атрибутной транслирующей грамматики

Атрибутная транслирующая грамматика - это транслирующая грамматика, обладающая следующими дополнительными свойствами:

1. Каждый символ грамматики (входной, нетерминальный и операционной) имеет конечное множество атрибутов, и каждый атрибут имеет (возможно, бесконечное) множество допустимых значений.
2. Все атрибуты нетерминальных и операционных символов делятся на синтезированные и унаследованные.
3. Унаследованные атрибуты подчиняются следующим правилам:
 - каждому вхождению унаследованного атрибута в правую часть правила вывода сопоставляется правило, позволяющее вычислить значение этого атрибута как функцию некоторых других атрибутов символов, входящих в левую или правую часть данного правила;
 - для каждого унаследованного атрибута начального символа грамматики задаётся начальное значение.
4. Правила вычисления значений синтезированных атрибутов определяются следующим образом:
 - каждому вхождению синтезированного атрибута нетерминального символа в левую часть правила вывода сопоставляется правило, позволяющее вычислить значение этого атрибута как функцию некоторых других атрибутов символов, входящих в левую или правую часть данного правила;
 - каждому синтезированному атрибуту операционного символа сопоставляется правило, позволяющее вычислить значение этого атрибута как функцию некоторых других атрибутов данного символа.
 - каждому синтезированному атрибуту операционного символа сопоставляется правило, позволяющее вычислить значение этого атрибута как функцию некоторых других атрибутов данного символа.

Будем записывать атрибуты в виде индексов соответствующих символов АТ- грамматики, при этом Для каждого атрибута будем указывать, к какому классу этот атрибут относится. Например:

$X_{a,b}$ синтезированный α унаследованный b

Правила вычисления атрибутов будем записывать в виде операторов присваивания, левая часть которых - атрибут или список атрибутов, а правая часть - функция.

Рассмотрим АТ-грамматику G_A , описывающую перевод оператора присваивания некоторого языка программирования в цепочку тетрад с кодами операций: СЛОЖИТЬ, УМНОЖИТЬ, ПРИСВОИТЬ. Левой частью оператора присваивания является идентификатор, а правой частью - бесскобочное арифметическое выражение, выполняемое слева направо в порядке написания операций.

Входными символами грамматики являются символы множества $\{i, +, *, =\}$, где i - лексема, представляющая идентификатор. Каждый идентификатор имеет один атрибут, значением которого является указатель на соответствующий элемент таблицы идентификаторов.

Операционные символы $\{+\}$, $\{*\}$, $\{:=\}$ соответствуют кодам операций тетрад, получаемых на выходе. Символы $\{+\}$ и $\{*\}$ имеют по три унаследованных атрибута, значениями которых являются указатели на элементы таблицы, в которой хранятся значения левого операнда, правого операнда и результата соответственно. Операционный символ $\{:=\}$ имеет два унаследованных атрибута. Значение первого атрибута - указатель на элемент таблицы идентификаторов, соответствующий идентификатору из левой Части оператора присваивания, а значением второго атрибута является указатель на элемент таблицы, в котором хранится значение выражения из правой Части оператора присваивания.

Словарь нетерминальных символов состоит из символов S , E и R . Начальный символ S не имеет атрибутов, символ E имеет один синтезированный атрибут, значением которого является указатель на элемент таблицы, содержащий значение выражения, порожденного E , а символ R имеет два атрибута: унаследованный (первый) и синтезированный. Значением унаследованного атрибута является указатель на элемент таблицы, соответствующий промежуточному результату, предшествующему R . Этот промежуточный результат является также левым операндом первой операции, порожденной нетерминальным символом R , если такая операция имеет место. Значение синтезированного атрибута символа R - это указатель на элемент таблицы, соответствующий значению подвыражения, которое получается после присоединения цепочки, порожденной символом R , к цепочке, представляющей левый операнд.

Правила вывода АТ-Грамматики G_A с соответствующими правилами вычисления следующих атрибутов:

E_t синтезированный t

$R_{p,t}$ унаследованный p синтезированный t

Атрибуты операционных символов унаследованные.

- (1) $S \rightarrow i_{p_1} := E_{q_1} \{:=\} p_2, q_2$

$$p_2 \leftarrow p_1$$

$$q_2 \leftarrow q_1$$
- (2) $E_{t_2} \rightarrow i_{p_1} R_{p_2}, t_1$

$$p_2 \leftarrow p_1$$

$$t_2 \leftarrow t_1$$
- (3) $R_{p_1, t_2} \rightarrow + i_{q_1} \{+\} p_2, q_2, r_1 R_{p_2, t_1}$

$$r_1, r_2 \leftarrow \text{GETNEW}$$

$$p_2 \leftarrow p_1$$

$$q_2 \leftarrow q_1$$

$$t_2 \leftarrow t_1$$
- (4) $R_{p_1, t_2} \rightarrow * i_{q_1} \{*\} p_2, q_2, r_1 R_{p_2, t_1}$

$$r_1, r_2 \leftarrow \text{GETNEW}$$

$$p_2 \leftarrow p_1$$

$$q_2 \leftarrow q_1$$

$$t_2 \leftarrow t_1$$
- (5) $R_{p_1, t_2} \rightarrow \epsilon$

$$p_2 \leftarrow p_1$$

В атрибутных правилах, связанных с правилами вывода (3) и (4), используется процедура-функция без параметров GETNEW , которая выдаёт значение указателя на свободную позицию таблицы, где будут запоминаться промежуточные результаты. Строго говоря, эта процедура не является функцией, т. к. при разных обращениях к ней получаются разные результаты. Таким образом, пользуясь процедурой GETNEW , мы несколько отступаем от определения АТ-грамматики. Вместо того чтобы использовать процедуру GETNEW , можно ввести дополнительные атрибуты для запоминания адресов занятых позиций таблицы. Однако при использовании процедуры GETNEW атрибутные правила получаются более простыми, и такой подход можно порекомендовать при практическом конструировании компиляторов.

АТ-Грамматики используются для получения атрибутных деревьев вывода, атрибутных активных цепочек и атрибутных переводов.

- 1) Процедура построения атрибутного дерева вывода включает следующие действия:
- 2) По соответствующей транслирующей грамматике построить дерево вывода активной цепочки, состоящей из входных и операционных символов без атрибутов.
- 3) Присвоить начальные значения унаследованным атрибутам начального символа грамматики.
- 4) Присвоить значения атрибутам входных символов, входящих в дерево вывода.
- 5) Найти атрибут, отсутствующий в дереве вывода, аргументы правила вычисления которого уже известны. Вычислить значение этого атрибута и добавить его в дерево.

Повторять шаг 4 до тех пор, пока значения всех атрибутов символов, входящих в дерево вывода, не будут вычислены или пока шаг 4 может выполняться.

На рис. 5.6 приведено атрибутное дерево вывода в АТ-грамматике, соответствующее входной цепочке $i_5 = i_7 + i_2 * i_3$ при условии, что процедура GETNEW выдаёт последовательные адреса ячеек, начиная с адреса 100. Заметим, что при применении правила (5) синтезированному атрибуту не терминала R присваивается значение его унаследованного атрибута. Затем это значение передаётся вверх по дереву вывода как синтезированный атрибут не терминала из левой части правил (4), (3) и (2) и используется в качестве второго атрибута операционного символа {ПРИСВОИТЬ}.

Последовательность атрибутных входных и операционных символов, полученная по атрибутному дереву вывода в АТ-грамматике, называется атрибутной активной цепочкой.

Множество пар, первым элементом которых является атрибутная входная цепочка, а вторым элементом — атрибутная последовательность операционных символов атрибутной активной цепочки, называется атрибутным

Пример 5.8

Грамматика GA для входной Цепочки $i_5 = i_7 + i_2 * i_3$ определяет следующую атрибутную последовательность операционных символов:

$$\{+\}_{7, 2, 100} \{*\}_{100, 3, 101} \{:=\}_{5, 101}$$

Определение: дерево вывода в АТ-грамматике называется завершённым, если в результате применения процедуры построения дерева значения всех атрибутов всех символов окажутся вычисленными.

Хотя для каждого атрибута, входящего в Дерево вывода, существует правило его вычисления, может оказаться, что между атрибутами возникла круговая зависимость, в результате чего нельзя получить завершённое дерево.

Определение: АТ-грамматика называется корректной тогда и только тогда, когда в результате применения описанной процедуры построения атрибутного дерева вывода получается завершённое дерево.

Далее будут рассмотрены два подкласса корректных АТ-Грамматик, которые часто используются при проектировании компиляторов: L-атрибутные транслирующие грамматики и S-атрибутные транслирующие Грамматики.

5.4.4 Вычисление значений атрибутов

Вернёмся ещё раз к проблеме вычисления значений атрибутов. После построения дерева вывода входной цепочки возникает вопрос о порядке вычисления значений атрибутов. По определению, атрибут можно вычислить, если известны значения всех атрибутов, от которых зависит его значение. Число деревьев вывода так же, как и входных цепочек, бесконечно, поэтому важно по самой грамматике уметь определять, является ли множество её правил вычисления атрибутов корректным.

Алгоритмы проверки корректности АТ-Грамматики приведены в [2, 25]. Они основаны на построении графа зависимостей и его анализе.

Узлами графа зависимостей служат атрибуты, которые нужно вычислить, а дугам ставятся в соответствие зависимости, определяющие, какие атрибуты вычисляются раньше, а какие позже.

Граф зависимостей R(D) представляет собой ориентированный Граф, который строится для некоторого дерева вывода D следующим образом:

- узлами R(D) являются пары (X, a) , где X - узел дерева D, а a - атрибут символа, служащего меткой узла X;
- дуга из узла (X_1, a_1) в узел (X_2, a_2) проводится, если семантическое правило, вычисляющее значение атрибута a_1 , непосредственно использует значение атрибута a_2 .

Вычисление значения двоичного числа по его символьному представлению можно задать при помощи следующей атрибутной грамматики [24]:

$$N_{v_4} \rightarrow L_{v_2, l_2, S_2} \cdot L_{v_3, l_3, S_{23}}$$

$$V_4 \leftarrow V_2 + V_3$$

$$S_2 \leftarrow 0$$

$$S_3 \leftarrow -l_3$$

$$N_{v_4} \rightarrow L_{v_2, l_2, S_2}$$

$$V_4 \leftarrow V_2$$

$$S_2 \leftarrow 0$$

$$L_{v_2, l_2, S_2} \rightarrow L_{v_3, l_3, S_3} \quad B_{v_1, S_1}$$

$$V_2 \leftarrow V_3 + V_1$$

$$S_1 \leftarrow S_2$$

$$S_3 \leftarrow S_2 + 1$$

$$l_2 \leftarrow l_3 + 1$$

$$L_{v_2, l_2, S_2} \rightarrow B_{v_1, S_1}$$

$$V_2 \leftarrow V_1$$

$$S_1 \leftarrow S_2$$

$$l_2 \leftarrow 1$$

$$B_{v_1, S_1} \rightarrow 0$$

$$V_1 \leftarrow 0$$

$$B_{v_1, S_1} \rightarrow 1$$

$$V_1 \leftarrow 1$$

В этой грамматике атрибуты имеют следующую семантику: v – значение (рациональное число), s – масштаб (целое число) и $-$ длина числа (целое число). Используя грамматику, построим дерево D вывода цепочки 101.01

Для построения узлов графа зависимостей необходимо последовательно рассмотреть вершины дерева D и для каждой из них построить столько узлов, сколько атрибутов имеет символ, которым помечена эта вершина.

Например, для корня дерева, обозначенного символом N с одним атрибуттом, нужно в граф включить один узел, поместив его парой (N, V_4) .

Для определения дуг графа необходимо использовать семантические правила. Например, рассмотрим узел Дерева зависимостей (N, V_4) . При построении дерева D непосредственные потомки корня N определялись правилом грамматики $N_{V_4} \rightarrow L_{V_2, l_2, S_2} \cdot L_{V_3, l_3, S_3}$. Атрибут V_4 не терминала N зависит, в соответствии с семантическим правилом $V_4 \leftarrow V_2 + V_3$, от атрибута V_2 левого сына, помеченного символом L , и атрибута V_3 правого сына. Поэтому в граф зависимостей необходимо включить дуги $((L, V_2), (N, V_4))$ и $((L, V_3), (N, V_4))$. Поступая аналогичным образом, мы можем построить граф зависимости.

В [24] Доказано, что атрибутная грамматика корректна, если граф зависимостей не имеет циклов.

Общие алгоритмы вычисления атрибутов весьма неэффективны, поэтому на практике используют методы, в которых для определения порядка вычисления не рассматриваются семантические правила. Например, если порядок вычисления атрибутов задаётся методом синтаксического анализа, то он не требует анализа семантических правил и построения графа зависимости.

5.5. Методика разработки описания перевода

При проектировании перевода следует руководствоваться следующими общими соображениями:

1. Описание перевода строится последовательно для конструкций входного языка в порядке их усложнения, начиная с простейших конструкций, например: выражение, оператор присваивания, оператор Цикла.

2. Выходная цепочка (результат перевода) включает в себя объекты (сущности) трех видов:

- объекты, передаваемые в выходную цепочку из входной цепочки;
- объекты, не изменяющиеся в процессе перевода (терминалы выходного языка);
- объекты, генерируемые во время перевода.

3. Транслирующая грамматика определяет порядок применения операционных символов (элементов перевода), строящих выходную цепочку, атрибуты же используются только для передачи значений символов из одних правил в другие.

4. Операционный символ, включаемый в правило вывода грамматики, может генерировать выходную цепочку до тех пор, пока в неё не должен быть включён объект, перемещаемый из входной цепочки и недоступный в данном правиле грамматики.

Проектирование описания перевода некоторой (одной!) конструкции входного языка выполняется в несколько этапов.

Неформальное описание перевода. Неформальное описание перевода представляет собой пару цепочек, первая из которых является конструкцией входного языка, а вторая – ее представлением в выходном языке. При разработке описания перевода рекомендуется:

- строить описание перевода на основе описания синтаксиса и семантики входного языка, учитывая все возможные форматы представления входной конструкции (например, if-then и if-then-else);
- конструкции, перевод которых уже описан и которые являются частью данной конструкции, считать терминальными;
- после построения неформального описания перевода выделить в нем символы, передаваемые из входной цепочки, и символы, генерируемые при его построении.

Описание синтаксиса. Описание синтаксиса выполняется в БНФ или расширенной БНФ. После стандартных преобразований БНФ преобразуется в КС-Грамматику, описывающую рассматриваемую конструкцию входного языка.

Определение транслирующей грамматики. При определении транслирующей грамматики требуется анализировать последовательность используемых при построении перевода правил грамматики, а также моделировать действия, выполняемые операционными символами грамматики. В связи с этим на данном этапе должны быть уточнены:

- метод синтаксического анализа, используемый при анализе и переводе Данной конструкции, и способ его реализации;
- интерпретация операционных символов (Цепочка символов, помещаемых на выходную ленту, или имя подпрограммы, которую необходимо выполнить).
- Исходными данными при построении транслирующей грамматики являются:
- КС-Грамматика, описывающая синтаксис конструкции;
- последовательность правил грамматики (разбор), используемых процессором при построении перевода;
- объекты, доступные процессору, выполняющему перевод (для процессора с магазинной памятью это текущий символ входной цепочки и, в зависимости от метода синтаксического анализа, весь магазин или его Часть).

Для построения транслирующей грамматики необходимо выполнить следующие действия:

1. Выбрать простейшую входную цепочку, соответствующую данной конструкции, и определить ее разбор.
2. В неформальном описании перевода выделить очередной (вначале первый) символ выходной цепочки, который должен быть передан в неё из входной цепочки.
3. На основании доступности (видимости) данных процессору и с учётом разбора, определить правило грамматики, в которое нужно включить элемент перевода (операционный символ), переносящий нужный символ в выходную цепочку. Если действия успешны, то перейти к шагу 5, иначе выполнить следующий шаг 4.
4. Если нельзя выполнить шаг 3, то соответствующий операционный символ включить в наиболее подходящее правило грамматики (в дальнейшем этот символ будет передан в него при помощи атрибутов), определив, если нужно, дополнительную выходную ленту процессора.
5. Возложить на новый (включённый) операционный символ действия по записи в выходную цепочку всех символов, вплоть до следующего символа, который должен быть передан из входной цепочки.
6. Повторить шаги 2-5 для всех символов, передаваемых из входной цепочки.
7. Тестирующую грамматику на более сложных, например, вложенных входных цепочках.

В простых случаях (простая конструкция входного языка, простое ее представление в выходном языке, удачные метод синтаксического анализа и его реализация, удачное расположение операционных символов) может оказаться, что для описания перевода достаточно только транслирующей грамматики.

Определение атрибутной транслирующей грамматики. При необходимости передачи данных между узлами синтаксического дерева транслирующую грамматику следует расширить до атрибутной транслирующей грамматики.

По определению, передача значений атрибутов в атрибутных транслирующих грамматиках возможна только в следующих направлениях:

- от символа левой части символам правой части правила вывода;
- от символов правой части правила символу левой части этого же правила вывода;
- между символами правой части правила вывода.

Для реализации передачи значения символа необходимо выполнить следующие действия:

1. Определить источник передаваемого значения (символ грамматики) и присвоить ему атрибут.
2. Определить приёмник значения (символ грамматики) и присвоить ему атрибут.
3. Учитывая ограничения в передаче атрибутов, определить маршрут передачи значения, для чего удобно пользоваться синтаксическим деревом грамматики.
4. Описать передачу значения атрибута правилами АТ-грамматики, стараясь, чтобы функции вычисления атрибутов были простейшими (копирующие правила).

Тестирование АТ-грамматики. Тестирование АТ-грамматики заключается в моделировании работы построенного по ней процессора, выполняющего перевод. При большом числе тестов и/или большой их длине эта работа может быть выполнена только при использовании соответствующих средств автоматизации. Для простых языковых конструкций, когда длина теста и их число невелики, а тестирование выполняется отдельно для каждой конструкции, начиная с простейших, работа эта не только необходима, но и реально выполнима.

Тестирование АТ-Грамматики позволяет определять два вида ошибок:

- ошибки в структуре выходной цепочки (неверный порядок следования символов в выходной цепочке);
- ошибки в значениях символов выходной цепочки, что связано с ошибками в передаче значений атрибутов.

При тестировании реального описания перевода желательно использовать комплексные тесты, определяющие оба вида ошибок.

5.6. Пример разработки АТ-грамматики

Выполним разработку описания перевода оператора цикла, имеющего следующий формат:

```
for <параметр> = <начальное значение> to <конечное значение> step <шаг>
<тело цикла> next <параметр>
```

Для упрощения введём следующие ограничения:

- <параметр> - идентификатор целой переменной;
- <начальное значение>, <конечное значение>, <шаг> - целые положительные константы;
- <тело цикла> - оператор цикла или терминал (не подлежащий переводу оператор);
- кодирование символов входной и выходной программы не рассматривается.

Оператор цикла, который требуется перевести, имеет вид:

```
for i = c1 to c2 step c3 <тело Цикла> next i
```

Допустим, что выходной язык не содержит оператора цикла. Тогда результат перевода (рис. 5.9) можно представить следующей последовательностью операторов:

```
i := c1; m1: if i > c2 then goto m2; <тело цикла>; i := i + c3; goto m1; m2:
```

Анализ выполняемого перевода позволяет сделать вывод о том, что выходная цепочка включает в себя объекты (сущности) трех видов:

1. Объекты, передаваемые в выходную цепочку из входной цепочки (переменная цикла i, начальное значение c1, конечное значение c2, шаг c3).

2. Объекты, не изменяющиеся в процессе перевода (множество терминальных символов выходного языка { ;, :, :=, +, >, if, then, goto }.

3. Объекты, генерируемые во время перевода (метки m1 и m2).

Составим БНФ, которая описывает синтаксис заданного оператора цикла, выбирая такую структуру описания, которая позволяет проиллюстрировать практически все проблемы синтеза АТ-грамматик:

- (1) <оператор цикла> ::= <заголовок цикла><тело цикла><конец цикла>
- (2) <заголовок цикла> ::= for <начальное значение><конечное значение><шаг>
- (3) <начальное значение> ::= <идентификатор> = <константа>
- (4) <конечное значение> ::= to <константа>
- (5) <шаг> ::= step <константа>
- (6) <тело цикла> ::= <оператор цикла>
- (7) <конец цикла> ::= next <идентификатор>

Введём обозначения: s - <оператор цикла>, A - <заголовок цикла>, B - <тело цикла>, C - <конец цикла>, D - <начальное значение>, E - <конечное значение>, F - <шаг>.

Тогда правила вывода КС-грамматики, описывающей синтаксис входного языка, имеют вид:

- (1) S → A B C
- (2) A → for D E F
- (3) D → id = cns
- (4) E → to cns
- (5) F → step cns
- (6) B → s
- (7) B → term
- (8) C → next id

При построении транслирующей грамматики требуется сделать некоторые предположения о ходе выполнения перевода и порядке анализа входной цепочки и построения выходной цепочки.

Для определённости будем считать, что:

используется восходящий метод синтаксического анализа, выполняемый процессором с магазинной памятью;

выполняется цепочечный перевод, в котором операционный символ вида {OUT := "str"} означает запись в выходную ленту OUT цепочки символов "str".

Определим действия, выполняемые первым по порядку элементом перевода (операционным символом), и правило, в которое он должен быть помещён. Для этого ещё раз рассмотрим неформальное описание перевода (см. рис. 5.9). Очевидно, что начать запись в выходную ленту процессор сможет только тогда, когда ему будет доступен параметр цикла id. Когда процессор выполняет свёртку, используя правило грамматики с номером (3), то id входит в основу и Доступен процессору.

Следовательно, в это правило мы можем включить операционный символ, который записывает на выходную ленту начальный участок перевода, вплоть до символа cns. Правило транслирующей грамматики с номером (3) будет выглядеть следующим образом:

- (3) D → id = cns { OUT:= "id := cns; m1 : if id >" }.

При выполнении свёртки по этому правилу операционный символ записывает на выходную ленту параметр цикла id. Запись этой переменной в OUT может быть выполнена функцией "читать третий сверху символ магазина процессора и записать его на выходную ленту". Аналогичные действия ("Читать верхний символ магазина и записать его на выходную ленту") выполняются для начального значения параметра цикла cns. Остальные символы не зависят от входной цепочки (терминалы выходного языка) и просто помещаются на выходную ленту.

Рассуждая подобным образом, мы получим следующие правила транслирующей грамматики, содержащие элементы перевода (в порядке их использования при анализе входной цепочки и построении выходной Цепочки):

- (4) E → to cns { OUT:= "cns then goto m2;" }
- (7) B → term { OUT:= "term;" }
- (5) F → step cns { OUT:= "id := id + cns; goto m1; m2 : " }

Добавив остальные правила грамматики, получим следующую транслирующую грамматику:

- (1) S → A B C
- (2) A → for D E F
- (3) D → id = cns { OUT:= "id := cns; m1 : if id >" }
- (4) E → to cns { OUT:= "cns then goto m2;" }
- (5) F → step cns { OUT:= "id := id + cns; goto m1; m2 : " }
- (6) B → S
- (7) B → term { OUT:= "term;" }
- (8) C → next id

При восходящих методах синтаксического анализа (см. гл. 8, 9) строится обращённый правый вывод (разбор) цепочки - последовательность номеров правил, используемых при свёртках, - который для рассматриваемого примера равен: ((3), (4), (5), (2), (7), (8), (1)). При этом на выходной ленте формируется цепочка:

id := cns1; m1: if id > cns2 then goto m2; id := id + cns3; goto m1; m2: term;

Анализ полученной цепочки позволяет сделать следующие выводы:

1. Тело цикла (терминальный оператор term) включается неправильно (не перед увеличением параметра цикла на значение шага, а в конец выходной цепочки). Это связано с тем, что правило (5) применяется раньше, чем правило (7).
2. При использовании правила (5) значение параметра цикла id недоступно (id вытолкнут из магазина при свёртке по третьему правилу грамматики). Значение id должно быть передано в элемент перевода правила (5) для правильной генерации выходной цепочки.
3. При переводе вложенных циклов возникнет конфликт меток m1 и m2, т. к. m1 и m2 - фиксированные значения. Значения меток должны генерироваться при выполнении перевода и передаваться от места возникновения метки к месту её использования.

Ошибки, обнаруженные в выходной цепочке, можно легко устранить, расширив транслирующую грамматику до атрибутной транслирующей грамматики. Первую ошибку можно исправить несколькими способами:

- включить в правило (4) атрибут, в который записать координаты выходной цепочки, куда должно быть вставлено тело цикла, передать значение этого атрибута в правило (5) и использовать его для выполнения включения тела цикла в выходную Цепочку;
- связать с не терминалом F атрибут, значением которого является строка символов. При использовании правила (5) записать в эту строку операторы, реализующие изменение параметра цикла, затем передать значение этого атрибута в правило (1), где и переписать его значение в выходную цепочку;
- пополнить атрибутный процессор, выполняющий перевод, Дополнительной лентой, на которую записывать операторы, реализующие изменение переменной Цикла в правиле (5). В конце перевода при свёртке по правилу (1) сцепить вспомогательную ленту с выходной лентой.

Применение первых двух способов не вызывает принципиальных затруднений, но неэффективно при реализации процессора, поэтому используем последний способ устранения ошибки: на выходную ленту OUT1 будем за- писывать операторы установки начального значения параметра цикла, проверки завершения цикла и тело цикла, а на дополнительную ленту OUT2 - операторы изменения переменной цикла с заданным шагом.

Вторая и третья ошибки устраняются стандартным способом: включением в правила грамматики атрибутов и функций их вычисления.

Преобразуем транслирующую грамматику в транслирующую атрибутную грамматику.

Рассмотрим правило транслирующей грамматики:

$$D \rightarrow id = cns \quad \{ OUT1 := "id := cns; m1 : if id >" \}$$

При выполнении свёртки по этому правилу значения терминальных символов id и cns требуется передать из входной цепочки в выходную. Это реализуется передачей соответствующих атрибутов в операционный символ и ис- пользованием их при построении выходной цепочки:

$$D \rightarrow id_{a1} = cns_{a2} \quad \{ OUT1 := "id_{n1} := cns_{n2}; m1: if id >" \}_{n1, n2} \\ s1, n1 \leftarrow a2; \quad n2 \leftarrow a2$$

Запись "id l" в операционном символе означает, что в данное место выходной цепочки нужно поместить идентификатор, значение которого определяется унаследованным атрибутом n1 (значения id и cns - это не числовые значения этих объектов, а ссылки на таблицы, которые содержат необходимую информацию).

Так как переменная цикла id используется и в других правилах грамматики (правила (5) и (8)), то ее значение нужно передать в эти правила, связав с символом D атрибут и определив функцию для его вычисления. Окончательно правило грамматики примет вид:

$$(3) D_{s1} \rightarrow id_{a1} = cns_{a2} \quad \{ OUT1 := "id_{a1} := cns_{n2}; m1: if id >" \}_{n1, n2} \\ s1, n1 \leftarrow a1; \quad n2 \leftarrow a2$$

Для определения способа передачи значения переменной цикла id (атрибут a1) в правила (5) и (8) грамматики, рассмотрим фрагмент дерева грамматики, изображённый на рис. 5.5.

Процесс передачи значения лексемы id (атрибут a1) из правила (3) грамматики отмечен на рис. 5.5 сплошными линиями:

1. Из правой части правила (3) значение атрибута входного символа a1 присваивается синтезированному атрибуту s1 символа D из левой части этого же правила.
2. Из правой части правила (2) значение синтезированного атрибута s1 нетерминального символа D присваивается унаследованному атрибуту n1 не терминала F из правой части этого же правила.
3. Из левой части правила (5) значение унаследованного атрибута n1 не терминала F передаётся в правую часть унаследованному атрибуту n2 операционного символа.

Соответствующие описанному процессу правила АТ-Грамматики имеют вид:

$$(3) D_{s1} = id_{a1} \quad \{ OUT1 := "id_{n1} := cns_{n2}; m1: if id >" \}_{n1, n2} \\ s1, n1 \leftarrow a1; \quad n2 \leftarrow a2$$

$$(2) A \rightarrow for D_{s1} E F_{n1} \\ n1 \leftarrow s1$$

$$(5) F_{n1} \rightarrow step cns \quad \{ OUT2 := " id := id +cns; goto m1; m2 : " \}_{n2}$$

$n2 \leftarrow n1$

Передача значения параметра цикла в правило (8) грамматики изображено на рис. 5.5 штриховыми линиями и описывается следующими правилами АТ-грамматики:

- (2) $A_{s2} \rightarrow \text{for } D_{s1} E F_{n1}$
 $s2, n1 \leftarrow s1$
- (1) $S \rightarrow A_{s2} B C_{n1}$
 $n1 \leftarrow s2$
- (8) $C_{n1} \rightarrow \text{next id } \{ \text{if } n2 \neq n3 \text{ then Error} \}$
 $n2 \leftarrow n1$

Особенностью правила (8) построенной АТ-Грамматики является то, что его операционный символ должен интерпретироваться не как элемент цепочечного перевода, а как выполнение действий, напрямую не связанных с генерацией выходной цепочки. В этот операционный символ передаются значения параметра цикла из заголовка цикла и оператора его завершения. Если эти значения не равны, то должно возникать состояние ошибки процессора (Error).

Процесс генерации и передачи меток иллюстрируется на рис. 5.12. Метки m1 и m2 генерируются в операционных символах правил (3) и (4) грамматики соответственно с помощью процедуры-функции NewLabel, возвращающей при обращении к ней уникальное значение метки. Правила АТ-грамматики, описывающие процесс генерации меток, имеют вид:

- (3) $D_{s1}, s2 \rightarrow id_{a1} = cns_{a2}$
{ OUT1 := "id_{n1} := cns_{n2}; Label_{s3}: if id >" }_{n1, n2, n3}
 $s1, n1 \leftarrow 01; n2 \leftarrow a2; s3 \leftarrow \text{NewLabel}; n3 \leftarrow s3; s2 \leftarrow n3;$
- (4) $E_{s3} \rightarrow \text{to } cns \{ \text{OUT1} := "cns_{n1} \text{ then goto Label}_{s1}; " \} n2$
 $s1 \leftarrow \text{NewLabel}; n2 \leftarrow s1; s3 \leftarrow n2;$
- (2) $A \rightarrow \text{for } D_{s1,s2} E_{s3} F_{n1, n2, n3}$
 $n1 \leftarrow s1; n3 \leftarrow s3$
- (5) $F_{n1, n2, n3} \rightarrow \text{step } cns \{ \text{OUT2} := "id := id + cns; goto Label}_{n4}; \text{Label}_{n5}: " \} n4, n5$
 $n4 \leftarrow n2; n5 \leftarrow n3$

Сцепление выходных лент OUT1 и OUT2 должно быть выполнено один раз сразу после анализа и перевода всей входной цепочки. Для реализации сцепления выходных лент грамматику необходимо пополнить новым (нулевым) правилом, включив в него операционный символ, выполнявши действия по конкатенации выходных лент. Результирующая АТ-грамматика будет иметь вид:

- (0) $S_0 \rightarrow S \{ \text{OUT1} := \text{OUT1} \parallel \text{OUT2} \}$
- (1) $S \rightarrow A_{s1} B C_{n1}$
 $n1 \leftarrow S1$
- (2) $A_{s4} \rightarrow \text{for } D_{s1,s2} E_{s3} F_{n1, n2, n3}$
 $s4, n1 \leftarrow s1; n2 \leftarrow s2; n3 \leftarrow s3;$
- (3) $D_{s1,s2} \rightarrow id_{a1} = cns_{a2}$
{ OUT1 := "id_{n1} := cns_{n2}; Label_{s3}: if id >" }_{n1, n2, n3}
 $s1, n1 \leftarrow 01; n2 \leftarrow 02; s3 \leftarrow \text{NewLabel}; n3 \leftarrow s3; s2 \leftarrow n3;$
- (4) $E_{s3} \rightarrow \text{to } cns_{a2}; \{ \text{OUT1} := "cns_{n1} \text{ then goto Label}_{s1}; " \} n1, n2$
 $n1 \leftarrow 01; s1 \leftarrow \text{NewLabel}; n2 \leftarrow s1; s3 \leftarrow n2;$
- (5) $F_{n1, n2, n3} \rightarrow \text{step } cns_{a1}$
{ OUT2 := "id_{n6} := id_{n6} + cns_{n7}; goto Label_{n5}; Label_{n4}: " }_{n4, n5, n6}
 $n4 \leftarrow n2; n5 \leftarrow n3; n6 \leftarrow n1; n7 \leftarrow a1;$
- (6) $B \rightarrow S$
- (7) $B \rightarrow \text{term}; \{ \text{OUT1} := "term"; \}$
- (8) $C_{n1} \rightarrow \text{next id}_{a1} \{ \text{if } n2 \neq n3 \text{ then Error} \}_{n4, n3}$
 $n2 \leftarrow n1; n3 \leftarrow a1$

При переводе вложенных операторов цикла использовании двух выходных лент в процессоре и их сцепление при выполнении перевода может привести к ошибкам в структуре выходной цепочки. рассмотрим перевод следующей входной цепочки:

```
for i1 = c5 to c12 step c13
    for i2 = c21 to c22 step c23
        top
```

next i2

next i1

разбор которой в виде последовательности номеров правил и синтаксического дерева приведён на рис. 5.13.

К моменту использования правила (0) на выходных лентах OUT1 и OUT2 будут находиться следующие цепочки (на данном этапе проверки считаем, что вычисление атрибутов выполняется правильно):

OUT1:

```
i1 := c5;
Label5: if i1 > c12 then goto Label12;
i2 := c21;
```

```

Label21: if i1 > c22 then goto Label22; top;
OUT2:
    i1 := i1 + c13; goto Label5; Label12:
    i2 := i2 + c23; goto Label21; Label22:
После выполнения сцепления выходная (основная) лента будет содержать перевод:
OUT1:
    i1 := c5;
Label5: if i1 > c12 then goto Label12;
    i2 := c21;
Label21: if i1 > c22 then goto Label22; top;
OUT2:
    i1 := i1 + c13; goto Label5; Label12:
    i2 := i2 + c23; goto Label21; Label22:
    i2 := c21;
Label21: if i1 > c22 then goto Label22;
    top;
    i1 := i1 + c13; goto Label5;
Label12:
    i2 := i2 + c23;
    goto Label21;
Label22:

```

Анализ выходной цепочки показывает, что перевод построен неверно: операторы завершения, внешнего и вложенного циклов следуют в обратном порядке. Связано это с тем, что вначале на выходную ленту OUT2 записываются операторы завершения внешнего цикла, и только после этого вложенного. При формировании же перевода на выходную ленту вначале должны быть записаны операторы завершения вложенного цикла, а затем - внешнего. Для того чтобы исправить эту ошибку, необходимо запись на ленту OUT2 и чтение из неё выполнять с одного конца, т. е. реализовать ею как стек с операциями PushOUT2 ("Втолкнуть на ленту OUT2") и PopOUT2 ("Вытолкнуть с ленты OUT2").

После включения операций PushOUT2 и PopOUT2 в правила (O) и (5) окончательно получим следующий вид атрибутной транслирующей грамматики:

- (0) $S_0 \rightarrow S \{ \text{while (OUT2 not empty)} \text{ OUT1 := OUT1 5 PushOUT2 } \}$
- (1) $S \rightarrow A_{s1} \ B \ C_{n1}$
 $n1 \leftarrow s1$
- (2) $A_{s4} \rightarrow \text{for } D_{s1,s2} \ E_{s3} \ F_{n1, n2, n3}$
 $s4, n1 \leftarrow s1; \quad n2 \leftarrow s2; \quad n3 \leftarrow s3;$
- (3) $D_{s1,s2} \rightarrow id_{n1} = cns_{a1}$
 $\{ \text{OUT1 := "id}_{n1} := cns_{n1}; \text{Label}_{s3}: \text{if id >"} \}_{n1, n2, n3}$
 $s1, n1 \leftarrow a1; \quad n2 \leftarrow a2; \quad s3 \leftarrow \text{NewLabel}; \quad n3 \leftarrow s3; \quad s2 \leftarrow n3;$
- (4) $E_{s3} \rightarrow \text{to } cns_{a1} \{ \text{OUT1 := "cnsnl then goto Label}_{s3}; \text{"} \}_{n1, n2}$
 $n1 \leftarrow a1; \quad s1 \leftarrow \text{NewLabel}; \quad n2 \leftarrow s1; \quad s3 \leftarrow n2;$
- (5) $F_{n1, n2, n3} \rightarrow \text{step } cns_{a1}$
 $\{ \text{PushOUT2("id := "id}_{n6} := id_{n6} + cns_{n7}; \text{goto Label}_{n5}; \text{Label}_{n4} :") \}_{n4, n5, n6}$
- (6) $B \rightarrow 8$
- (7) $B \rightarrow \text{term; } \{ \text{OUT1 := "term; } \}$
- (8) $C_{n1} \rightarrow \text{next } id_{a1} \{ \text{if } n2 \neq n3 \text{ then Error } \}_{n2, n3}$
 $n2 \leftarrow n1; \quad n3 \leftarrow a1$

Читателям рекомендуется самостоятельно проверить правильность окончательного варианта построения АТ-Грамматики для различных тестовых входных цепочек.

Контрольные вопросы

1. Как определяется понятие семантики языков программирования?
2. Как можно описать простейшие переводы?
3. Какие переводы называются синтаксически управляемыми?
4. Что такое элемент перевода?
5. Дайте определение СУ-схемы. Какие цепочки называются входными и выходными?
6. Как определить входную и выходную грамматики СУ-схемы?
7. Как происходит порождение пар цепочек под управлением СУ-схемы?
8. Как происходит преобразование деревьев под управлением СУ-схем?
9. Какие СУ-схемы называются простыми?
10. Какие переводы позволяют описать простые СУ-схемы?
11. Какая грамматика называется транслирующей гамматикой?
12. Что такое входная и выходная грамматики транслирующей грамматики?
13. Как определяется активная цепочка, ее входная и операционная части?
14. Как связаны вывод цепочки и дерево вывода активной цепочки?
15. Какая Грамматика называется атрибутной транслирующей грамматикой?
16. Какие атрибуты называются унаследованными и синтезированными?

17. В чем основные различия между унаследованными и синтезированными атрибутами?
18. Что представляет собой правило вычисления атрибутов?
19. Как строится граф зависимостей?
20. Как график зависимостей используется при вычислении атрибутов?
21. Какие виды объектов включаются в выходную цепочку при переводе?
22. Из каких этапов состоит процесс построения транслирующей грамматики?
23. Из каких этапов состоит процесс построения АТ-Грамматики?
24. Какие типы ошибок позволяет найти тестирование АТ-Грамматики?

Упражнения

1. Разработайте простую СУ-схему, описывающую перевод арифметических скобочных выражений, содержащих операции '+' и '*' в:
 - 1.1. постфиксную запись.
 - 1.2. префиксную запись.

2. В языке ALGOL-60 выражения строятся с помощью операций, имеющих следующие приоритеты (в порядке убывания):

(1) ↑	(4) ≤ < = ≠ > ≥	(7) √
(2) */ ÷	(5) ¬	(8) →
(3) + -	(6) ∧	(9) ≡

Разработайте простую СУ-схему, описывающую перевод инфиксных выражений, содержащих перечисленные операции, в постфиксную запись.

3. Постройте МП-преобразователи, реализующие СУ-схемы из упр. 1.
4. Для простой СУ-схемы:

$$E \rightarrow + EE, EE +$$

$$E \rightarrow * EE, EE *$$

$$E \rightarrow a, a$$

постройте эквивалентный МП-преобразователь.

5. Для заданной СУ-схемы:

$$S \rightarrow a S^{(1)} b A c S^{(2)} d S^{(3)}, e A c S^{(2)} d S^{(3)} f S^{(1)} g$$

$$S \rightarrow i, i$$

$$A \rightarrow i, i$$

постройте переводы для входных цепочек:

$$5.1. aibicd .$$

$$5.2. aaibicidibicid .$$

6. Постройте вывод в транслирующей грамматике G0T цепочек:

$$6.1. i^* (i+i).$$

$$6.2. (i+i)^* (i + i).$$

7. Определите транслирующую грамматику, допускающую в качестве входа произвольную цепочку из нулей и единиц и порождающую на выходе:

$$7.1. \text{ обращение входной цепочки.}$$

$$7.2. \text{ цепочку } 0n1m, \text{ где } n - \text{ число нулей, а } m - \text{ число единиц во входной цепочке.}$$

8. Постройте транслирующую грамматику, определяющую перевод логических выражений, составленных из логических переменных, скобок и знаков операций дизъюнкции, конъюнкции и отрицания:

$$8.1. \text{ из инфиксной записи в ПОЛИЗ.}$$

$$8.2. \text{ из ПОЛИЗ в инфиксную запись.}$$

8.3. из инфиксной записи в функциональную запись такую, например, что выражение $a \cap (b \cup c)$ будет представлено в виде $f\cap(a, f\cup(b, c))$, где $f\cap$ и $f\cup$ - отдельные символы.

9. Постройте АТ-Грамматику, описывающую перевод оператора присваивания некоторого

Гипотетического языка программирования в цепочку тетради с кодами операций: ПРИСВОИТЬ, СЛОЖИТЬ, ВЫЧЕСТЬ, УМНОЖИТЬ, ДЕЛИТЬ. Левой частью оператора присваивания является идентификатор, а правой частью бесскобочное арифметическое выражение, выполняемое справа налево в порядке написания операций. В арифметическом выражении можно использовать идентификаторы и знаки арифметических операций: +, -, *, /.

10. Постройте АТ-грамматику, описывающую перевод оператора присваивания некоторого гипотетического языка программирования в цепочку тетрад с кодами операций: ПРИСВОИТЬ, И, ИЛИ, НЕ. левой частью оператора присваивания является идентификатор, а правая часть представляет собой логическое выражение, составленное из идентификаторов, скобок и знаков логических операций: ∩, ∪, ¬. порядок выполнения логического выражения определяется обычным приоритетом логических операций.

11. Постройте АТ-грамматику, описывающую перевод в тетрадку с кодом операции ПЕРЕХОД_ПО_РАВНО условного оператора, которому соответствует следующее правило входной Грамматики:

условный оператор $\rightarrow \text{if} = \text{goto} .$

Лексема представляет собой номер оператора, к которому осуществляется переход в случае равенства, а - идентификатор.

12. Для входной грамматики, описывающей арифметические выражения, с правилами вывода:

$S \rightarrow E$ $E \rightarrow (E)$

$E \rightarrow E := E$ $E \rightarrow i$

$E \rightarrow E + E$

где - лексема, значением которой является указатель на элемент таблицы идентификаторов, а семантика выражений такая же, как в языке С, постройте АТ-грамматику, которая проверяет, представляет ли левая Часть выражения значение, и, в случае успеха, стройте ПОЛИЗ выражения.

13. Для входной грамматики, описывающей объявление переменных, с правилами вывода:

(1) $D \rightarrow L$ (4) $t \rightarrow \text{int}$

(2) $L \rightarrow , L$ (5) $t \rightarrow \text{real}$

(3) $L \rightarrow :$

постройте АТ-грамматику, которая заносит значение типа каждого идентификатора в таблицу идентификаторов.

14. Для входной грамматики, описывающей арифметические выражения, с правилами вывода:

$E \rightarrow (E + E)$

$E \rightarrow (E * E)$

$E \rightarrow i$

где - лексема, значением которой является указатель на элемент таблицы идентификаторов, определить постфиксную 8-атрибутную грамматику, описывающую перевод этих выражений в цепочку тетради с концами операций: СЛОЖИТЬ и УМНОЖИТЬ, и построить 8-атрибутный ДМП-процессор, выполняющий заданный перевод.

UFPHOBIAK

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

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

1. 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.
2. Аргументами правила вычисления значения синтезированного атрибута символа из левой части правила вывода являются унаследованные атрибуты этого символа или произвольные атрибуты символов из правой части.
3. 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^A .

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.

Определение: 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.

6.2 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 $x, y \leftarrow 1$.

Определение: 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 $a, b, x \leftarrow y$.

Определение: 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.

Определение: 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 $G4$.

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

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:

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

Для каждого некопирующего правила $z_1, z_2, \dots, z_b \leftarrow f(y_1, y_2, \dots, y_n)$ 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 $\{f\} a_n \leftarrow a_2, \dots, a_n$. Заменить некопирующее правило на следующие ($n + 1$) копирующих правил:

$a_i \leftarrow y_n$ 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.

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 синтезированный p
 $\{ОТВЕТ\}p$, унаследованный r
 $S \rightarrow E_p \{ОТВЕТ\}r$
 $r \leftarrow p$
 $E_p \rightarrow + E_q E_p$
 $p \leftarrow q + r$
 $E_p \rightarrow * E_q E_r$
 $p \leftarrow q * r$
 $E_p \rightarrow c_r$
 $p \leftarrow r$

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 $\{ОТВЕТ\}$ также имеют по одному атрибуту. Значением синтезированного атрибута символа E является значение подвыражения, порожденного этим символом, а значением унаследованного атрибута символа $\{ОТВЕТ\}$ — значение всего выражения, порожденного грамматикой.

Исходная грамматика содержит два некопирующих 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 $\{\text{СЛОЖИТЬ}\}A, B, R$ атрибутное правило имеет вид: $R \leftarrow A + B$, а для операционного символа $\{\text{УМНОЖИТЬ}\}A, B, R$ — $R \leftarrow A \times B$.

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 единственное возможное место расположения этого символа. Полученная в результате преобразования L-атрибутная грамматика в форме простого присваивания имеет вид:

E_p синтезированный p
 $\{ОТВЕТ\}r$, 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
 $S \rightarrow E_p \{ОТВЕТ\}r$
 $r \leftarrow p$
 $E_p \rightarrow + E_q E_p \{\text{СЛОЖИТЬ}\}A, B, R$
 $A \leftarrow q$
 $B \leftarrow r$
 $P \leftarrow R$
 $E_p \rightarrow * E_q E_p \{\text{УМНОЖИТЬ}\}A, B, R$
 $A \leftarrow q$

$B \leftarrow r$
 $p \leftarrow R$
 $E_p \rightarrow c_r$
 $r \leftarrow r$

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:

(β, y) , где $\beta \in V_p^*$ — цепочка из правой части правила транслирующей грамматики $A \rightarrow y\beta$, а $y \in \Sigma_a^*$;

ВЫДАЧА (X), где $X \in \Sigma_a$

ВЫБРОС;

ДОПУСК;

ОШИБКА.

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

- ($x, X\alpha, \pi \vdash (x, \beta\alpha)$, 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.
- ($x, X\alpha, \pi \vdash (x, \alpha, \pi X)$, если $M(X, a) = \text{ВЫДАЧА } \{X\}$). 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:

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.

$M(X, a) = \text{ВЫДАЧА}(X)$ для всех $x \in \Sigma_a$ И $a \in \Sigma_i \cup \{\epsilon\}$.

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

$M(\perp, \epsilon) = \text{допуск}$

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^1 путем включения в нее операционных символов $\{i\}$, $\{+\}$ и $\{* \}$:

$E \rightarrow TE' \quad (5) \quad T' \rightarrow * P \{ *\} T'$

$E' \rightarrow + T \{ + \} E' \quad (6) \quad T' \rightarrow \epsilon$

$$\begin{array}{ll} E' \rightarrow \epsilon & (7) \quad P \rightarrow i \{i\} \\ T \rightarrow PT' & (8) \quad P \rightarrow (E) \end{array}$$

Управляющая таблица должна содержать 14 строк, помеченных символами из множества $N \cup \sum_i \cup \Sigma_a \cup \{\perp\}$, и 6 столбцов, помеченных символами из множества $\sum_i \cup \{\epsilon\}$.

Построение управляющей таблицы для строк, отмеченных символами из множества $N \cup \sum_i \cup \{\perp\}$, ничем не отличается от построения таблицы для соответствующей входной LL(1)-грамматики (табл. 1.1), а строки управляющей таблицы, отмеченные операционными символами, содержат значения 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	*P{*} T', ϵ		ϵ, ϵ
T	PE', ϵ	PE', ϵ				
T'			ϵ, ϵ	ϵ, ϵ	*P{*} T', ϵ	ϵ, ϵ
P	i{i}, ϵ	(E), ϵ				
i	ВЫБР ОС					
(ВЫБР ОС				
)			ВЫБРОС			
+				ВЫБРОС		
*					ВЫБР ОС	
\perp						ДО ПУ СК
{i}	ВЫДАЧА({i})					
{+}	ВЫДАЧА({+})					
{*}	ВЫДАЧА({*})					
Начальное содержимое магазина - E \perp						

Начальное содержимое магазина — E \perp

Для входной цепочки $i + i * i$ ДМП-процессор проделает следующую последовательность тактов:

$$\begin{aligned} & (i + i * i, E\perp, \epsilon) \vdash (i + i * i, TE'\perp, \epsilon) \\ & \vdash (i + i * i, PT'E'\perp, \epsilon) \\ & \vdash (i + i * i, \{i\} TE'\perp, \epsilon) \\ & \vdash (+ i * i, \{i\} T'E'\perp, \epsilon) \\ & \vdash (+ i * i, \{i\} TE'\perp, \{i\}) \\ & \vdash (+ i * i, \{i\} E'\perp, \{i\}) \\ & \vdash (+ i * i, +T\{+\} E'\perp, \{i\}) \\ & \vdash (i * i, T\{+\} E'\perp, \{i\}) \\ & \vdash (i * i, PT'\{+\} E'\perp, \{i\}) \\ & \vdash (i * i, i\{i\} T'\{+\} E'\perp, \{i\}) \\ & \vdash (*i, \{i\} T'\{+\} E'\perp, \{i\}) \\ & \vdash (*i, T'\{+\} E'\perp, \{i\} \{i\}) \\ & \vdash (*i, *P{*}T'\{+\} E'\perp, \{i\} \{i\}) \\ & \vdash (i, P{*}T'\{+\} E'\perp, \{i\} \{i\}) \\ & \vdash (i, i\{i\} {*}T'\{+\} E'\perp, \{i\} \{i\}) \\ & \vdash (\epsilon, \{i\} {*}T'\{+\} E'\perp, \{i\} \{i\}) \\ & \vdash (\epsilon, \{i\} {*}T'\{+\} E'\perp, \{i\} \{i\}) \\ & \vdash (\epsilon, T'\{+\} E'\perp, \{i\} \{i\} \{i\} \{*\}) \\ & \vdash (\epsilon, \{+\} E'\perp, \{i\} \{i\} \{i\} \{*\}) \\ & \vdash (\epsilon, E'\perp, \{i\} \{i\} \{i\} \{i\} \{*\}) \\ & \vdash (\epsilon, \perp, \{i\} \{i\} \{i\} \{i\} \{*\}) \end{aligned}$$

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 $\text{FIRST}(x) = a$. Опишем действия, которые должен выполнять L-атрибутный ДМП-процессор в зависимости от элемента управляющей таблицы $M(X, o)$, 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.

$M(X, a) = (\beta, y)$ (в верхушке магазина находится нетерминал). В этом случае L-атрибутный ДМП-процессор вталкивает в магазин цепочку символов β из правой части распознаваемого правила $B \text{ DSLFTN } WTGJXRE$ операционных символов y . 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 нетерминалы унаследованные 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 $M(X, a) = (\beta, y)$.

При выполнении перехода L-атрибутный ДМП-процессор выполняет следующие атрибутные действия:

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.

Если источник копирующего правила доступен, то значение источника помещается в соответствующее поле магазинного символа.

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.

Таблица 1.5. Управляющая таблица для транслирующей грамматики, полученной из исходной грамматики

S	i:=E{				
E	iR, ε				
R		+	*i{	ε, ε	
/	ВЫБ				
:		В			
+			ВЫ		
*				ВЫ	
⊥					ДОПУС
{	ВЫДАЧА({:=}р q)				
{	ВЫДАЧА({+}Р, q. г)				
{	ВЫДАЧА({*}Р. q. г)				

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 G1 и имеет следующие правила:

$$\begin{aligned} S &\rightarrow E\perp \\ E &\rightarrow TE' \\ E' &\rightarrow +T\{+\} \quad E' \mid \epsilon \\ T &\rightarrow PT' \\ T' &\rightarrow *P\{*\} \quad T' \mid \epsilon \\ P &\rightarrow (E) \mid i \{i\} \end{aligned}$$

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})
end
else
Error
end; {Proc_P}
begin
{В начале анализа переменная Symb содержит первый символ цепочки}
Proc_E;
if Symb = 'I' then Access
else
Error
end {Recurse_Method_TG};

```

На рис. 1.2 приведен список имен процедур в порядке их вызова при переводе входной цепочки $(i + i) * i \perp$. в предпоследнем столбце в строке, соответствующей вызову процедуры NextSymb, изображена непрочитанная часть входной цепочки (текущий символ находится слева) непосредственно после вызова этой процедуры, а в последнем столбце в строке, соответствующей вызову процедуры Output(Y), изображена часть выходной цепочки, построенной непосредственно после вызова процедуры Output(Y).

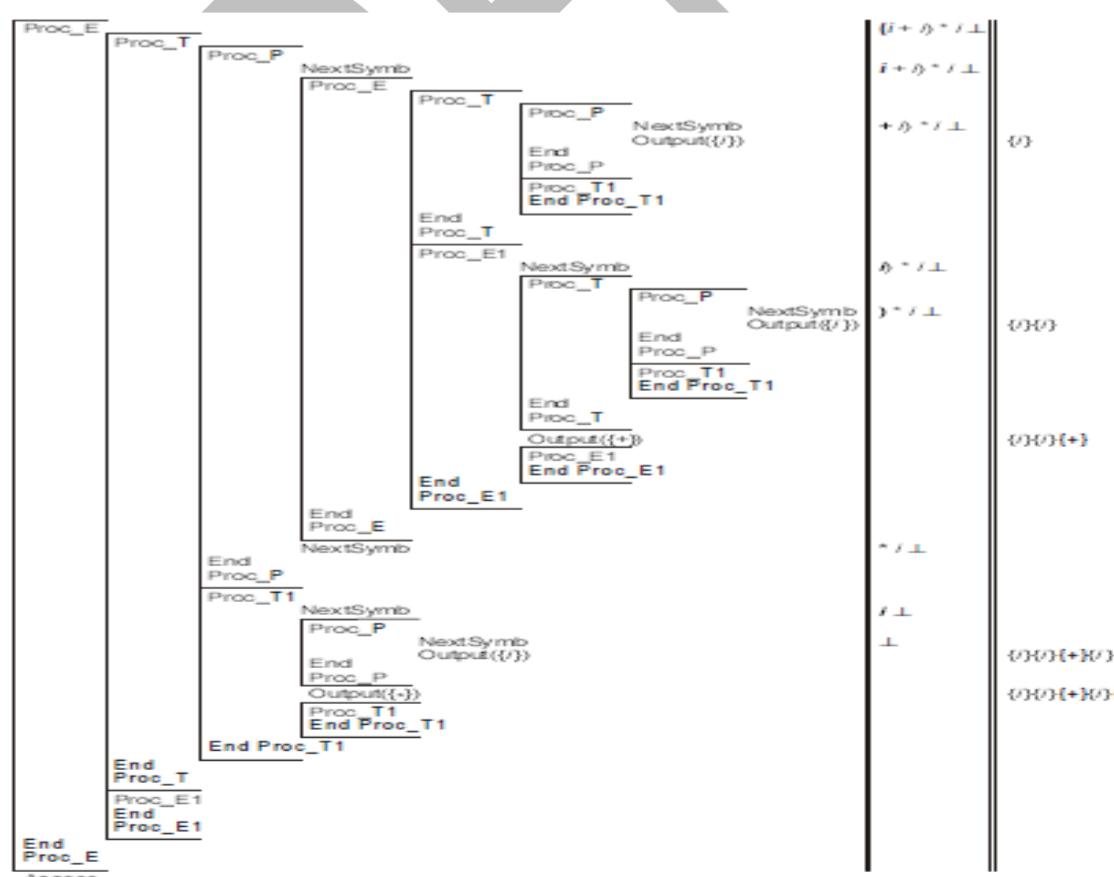


Рисунок 1.2. Порядок вызова процедур при переводе цепочки $(i + i) * i \perp$. Для того, чтобы методом рекурсивного спуска можно было выполнять L-атрибутный перевод, введем в процедуры распознавания нетерминальных символов параметр для каждого атрибута этого символа. При этом если параметру процедуры 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 нетерминалных 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:

$$\begin{aligned} R_{p_1, t_2} &\rightarrow + i_{q_1} \{*\}_{p_2, q_2, r_1} R_{r_1, t_1} \\ R_{p_1, t_2} &\rightarrow + i_{q_1} \{*\}_{p_2, q_2, r_1} R_{r_1, t_1} \\ R_{p_1, p_2} &\rightarrow \epsilon \end{aligned}$$

т. к. первые два вхождения нетерминала R имеют атрибуты p_1, t_2 , а последнее вхождение — p_1, p_2 . В этом случае необходимо выбрать какой-то один список имен атрибутов нетерминала R (например, р и t), использовать эти имена при описании типа атрибутов этого нетерминала (например, унаследованный р, 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.

$$\begin{aligned} A &\rightarrow B_x C_y D_z \\ y, z &\leftarrow x \end{aligned}$$

Атрибутам x, y и z присваивается одно и то же значение, поэтому им можно дать общее имя. Используя новое имя a, получим правило вывода грамматики ($A \rightarrow B_x C_y D_z$), которое не требует правила вычисления атрибута x.

$$\begin{aligned} A_x &\rightarrow B_y \{f\}_z \\ y, z &\leftarrow x \end{aligned}$$

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.

$$\begin{aligned} A_{x, y} &\rightarrow a, B_z C_t \\ y, z, t &\leftarrow x \end{aligned}$$

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:

$$\begin{aligned} A_{x, y} &\rightarrow a, B_y C_y \\ y &\leftarrow x \end{aligned}$$

Новый способ записи имен атрибутов позволяет непосредственно превращать списки атрибутов нетерминалных символов грамматики в списки параметров процедур для распознавания нетерминалных 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 вхождениям нетерминала в левые части правил вывода, становится списком формальных параметров соответствующей процедуры; 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 атрибутным 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 атрибутного правила (если соглашение об одинаковых именах выполнено, то в левой части атрибутного 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 атрибутного 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:

Et синтезированный t

Rp, 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) $R_{p,t} \rightarrow I_{q1} \{+\} p, q_1, r_1 R_{r1}, t$
 $r1 \leftarrow \text{GETNEW}$
- (4) $R_{p,t} \rightarrow X I_{q2} \{*\} p, q_1, , r_1 R_{r1}, t$
 $r2 \leftarrow \text{GETNEW}$
- (5) $R_{p,t} \rightarrow \epsilon$
 $t \leftarrow p$

С учетом выполненных преобразований процедура на языке Pascal, реализующая атрибутный 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
```

```

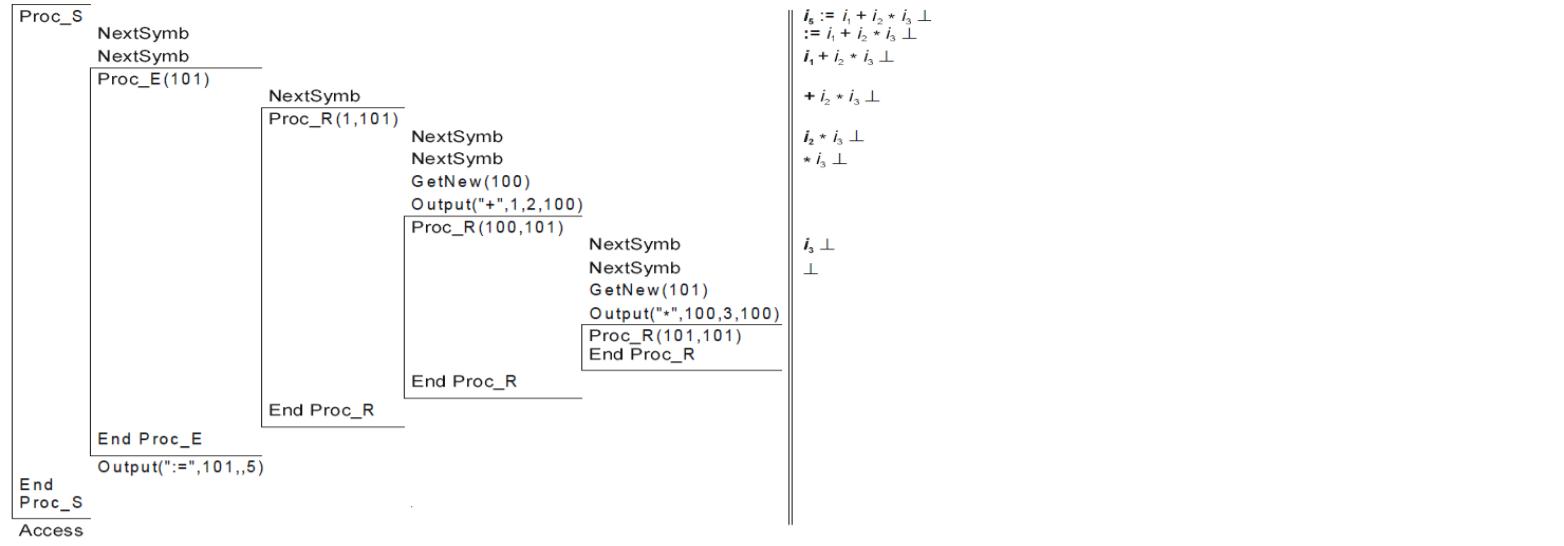
        NextSymb;
        Proc_E(b);
        Output({:=}, a, b)
    end
    else
        Error
    end
else
    Error
end; {Proc_S}
Procedure Proc_E (var c: integer);
    var d: integer; {локальная переменная Proc_E}
begin
    if ClassSymb = ClassId {текущий символ - идентификатор} then
        begin
            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 {Recurse_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.



Математическая модель восходящего ДМП-процессор

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 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 $G_q = (\{E, T \setminus P\}, \{i, +, *, (,)\}, P, E)$, где $P = \{E \rightarrow E + T, E \rightarrow T, T \rightarrow T^* P, P \rightarrow i, P \rightarrow (E)\}$, будет выглядеть следующим образом:

```

E → E + T{СВЕРТКА, 1}
E → T{СВЕРТКА, 2}
T → T* P {СВЕРТКА, 3}
T → P {СВЕРТКА, 4}
P → i {СВЕРТКА, 5}
P → (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, порожденная 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:

```

E → E+ T{+}
E → T
T → T*P {*}
T → P
P → i{i}
P → (E)

```

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.

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 p атрибутами представляется в магазине ($p + 1$)-ой ячейками, верхняя из которых содержит имя символа, а 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.

На рис. 6.6 изображена последовательность состояний магазина S-атрибутного ДМП-процессора, осуществляющего перевод Цепочки $i_2 + i_3 \times i_5$ в последовательность тетради со знаками операций СЛОЖИТЬ и УМНОЖИТЬ. S- атрибутная транслирующая Грамматика, входной грамматикой которой является основная грамматика, имеет вид:

E_p - СИНТЕЗИРОВАННЫЙ р

- (1) $E_{p_2} \rightarrow E_{q_1} + E_{t_1} \{+\}_{q_2, t_2, p_1}$
 $q_2 \leftarrow q_1$
 $t_2 \leftarrow t_1$
 $p_1, p_2 \leftarrow \text{GETNEW}$
- (3) $E_{p_2} \rightarrow E_{q_1} * E_{t_1} \{*\}_{q_2, t_2, p_1}$
 $q_2 \leftarrow q_1$
 $t_2 \leftarrow t_1$
 $p_1, p_2 \leftarrow \text{GETNEW}$
- (5) $E_{p_2} \rightarrow p_1$
 $p_2 \leftarrow p_1$
- (6) $E_{p_2} \rightarrow (E_{p_1})$
 $p_2 \leftarrow p_1$

На рис. 6.6, а приведено начальное состояние магазина. Рис. 6.6, б представляет результат операции ПЕРЕНОС входного символа с атрибутом 2 в магазин, а рис. 6.6, в результат выполнения операции СВЕРТКА, выполняемой в соответствии с правилом (5). Рис. 6.6, г и д иллюстрируют результат выполнения операции ПЕРЕНОС символов '+' и i_3 соответственно из входной строки в магазин, а рис. 6.6, е - результат выполнения операции СВЕРТКА символа i_3 в E_3 . Аналогично выполняются операции переноса и свёртки для символов '*' и i_5 (рис. 6.6, ж-е).

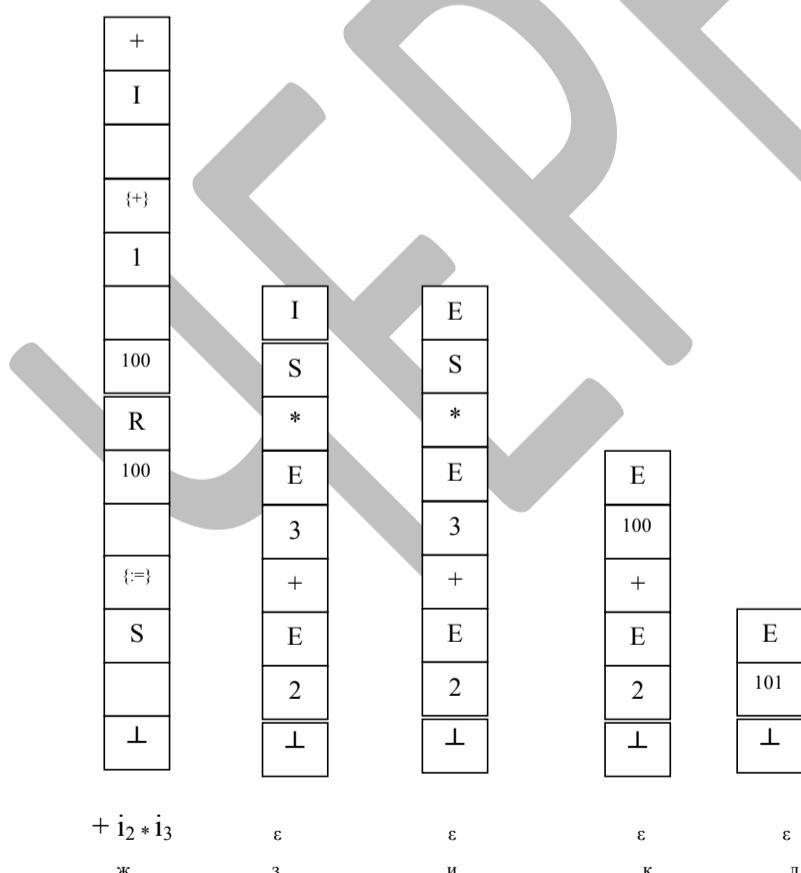


Рис. 6.6. Последовательность состояний магазина при обработке входной цепочки $i_2 + i_3 * i_5$

Рис. 6.6, и иллюстрирует результат выполнения операции СВЕРТКА, выполняемой в соответствии с правилом (1). При этом поля атрибутных символов из правой Части этого правила используются для формирования тетрады,

*	3	5	100
---	---	---	-----

которая выдаётся в выходную цепочку.

Далее выполняется операция СВЕРТКА по правилу (3) с выдачей тетрады:

+	2	100	101
---	---	-----	-----

(рис. 6.6, к-л).

Контрольные вопросы

1. Дайте определение L-атрибутной и S-атрибутной транслирующих грамматик.
2. Какое правило вычисления атрибутов называется копиющим?
3. В каком случае АТ-грамматика имеет форму простого присваивания? Приведите процедуру преобразования произвольной АТ-грамматики в форму простого присваивания.
4. Приведите алгоритм построения управляющей таблицы для транслирующей грамматики Цепочечного перевода, входной грамматикой которой является LL(!)-Грамматика.
5. Опишите, каким образом представляются в магазине и вычисляются унаследованные и синтезированные атрибуты символов Грамматики в L-атрибутном ДМП-процессоре.
6. Приведите процедуру преобразования ДМП-процессора в L-атрибутный ДМП-процессор.
7. Как программируются процедуры в методе рекурсивного спуска, реализующего перевод, описываемый транслирующей грамматикой?
8. С какой целью осуществляется переименование имён атрибутов нетерминальных символов из левых частей правил вывода L-атрибутной транслирующей грамматики при реализации вывода L-атрибутного перевода методом рекурсивного спуска?
9. Как программируются процедуры в методе рекурсивного спуска, реализующего перевод, описываемый L-атрибутной транслирующей Грамматикой?
10. Каким образом восходящий анализатор можно расширить действиями по выполнению перевода?
11. Опишите, каким образом представляются в магазине и вычисляются унаследованные и синтезированные атрибуты символов грамматики в S-атрибутном ДМП-процессоре.

Упражнения

1. Следующие правила вывода грамматики являются частью некоторой атрибутной грамматики. Атрибуты p, q и r – унаследованные, a, s и f – синтезированные. Для каждого правила определите, от каких атрибутов могут зависеть правила вычисления p и r , чтобы это правило было L-атрибутным?
 - 1.1. $As, q \rightarrow ai At, p Br.$
 - 1.2. $As, q \rightarrow \{c\}p Bi At, r.$
 - 1.3. $As, q \rightarrow si Br At, p.$
2. Приведите следующие правила вывода АТ-Грамматики к форме простого присваивания (имена унаследованных атрибутов начинаются с символа i , а имена синтезированных атрибутов – с символа s):
 - 2.1. $As_1, i_1 \rightarrow E$
 $s_1 \leftarrow \sin(i_1)$
 - 2.2. $E_{s_1}, i_1 \rightarrow As_2 Bs_3, i_2 C_{s_4}, i_3 D_{s_5}, i_4$
 $i_2 \leftarrow i_1$
 $i_3 \leftarrow i_1 * i_1$
 $i_4 \leftarrow i_2 * i_3$
 $s_1 \leftarrow s_2 * s_3$
3. Преобразуйте L-атрибутную грамматику цепочечного перевода к форме простого присваивания и постройте для нее L-атрибутный ДМП-процессор.

A_p СИНТЕЗИРОВАННЫЙ p
 B_p СИНТЕЗИРОВАННЫЙ p
 $\{b\}_p$ УНАСЛЕДОВАННЫЙ p

$(1) S \rightarrow a_p A_q B_r A_s \{b\}_t$ $r \leftarrow p + 2q$ $t \leftarrow r + s$	$(3) A_p \rightarrow B$ $p \leftarrow B$
--	---

$$(2) A_p \rightarrow a_q A_r B_s B_t B_u S$$
$$s \leftarrow q + r$$
$$t \leftarrow q * s - 4$$
$$u, p \leftarrow r + t$$

$$(4) B_p \rightarrow a_q A_r \{b\}_t$$
$$t \leftarrow r + q$$

4. Для Грамматики, заданной в упр. 2, постройте процессор методом рекурсивного спуска.

5. Для входной грамматики, описывающей синтаксис арифметических выражений, с правилами вывода

$$E \rightarrow (E + E)$$
$$E \rightarrow (E * E)$$
$$E \rightarrow i$$

Где i – лексема, значением которой является указатель на элемент таблицы идентификаторов, определите постфиксную S -атрибутную грамматику, описывающую перевод этих выражений в цепочку тетрад с кодами операций: СЛОЖИТЬ и УМНОЖИТЬ, и постройте S -атрибутный ДМП-процессор, выполняющий заданный перевод.

ЧЕРНОВЫЙ

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.WriteLine("(" + this.right[i]);
        else
            System.Console.Write(", " + this.right[i]);
    }
    if (i == 0)
        System.Console.WriteLine();
    else
        System.Console.WriteLine(")");
    }
}

```

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 fname) {
        System.Console.WriteLine(fname);
        // 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
    }
}

```

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(n2^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); // булеван
}
} // 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.Write("      (" + 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..

- A). УстраниТЬ из грамматики G бесполезные символы. Применить алгоритм 3.3. к грамматике G.
- B). УстраниТЬ из грамматика G ε-правила, применить алгоритм 3.5.
- C). УстраниТЬ из КС грамматики G цепные правила, применить алгоритм 3.6.
- D). УстраниТЬ левую рекурсию в заданной КС-грамматике G_1 , порождающей скобочные арифметические выражения. Применить алгоритм 3.7. к грамматике G.
- F). Определить в какой форме (Грейбах, Хомского) находится КС-грамматика G' .
- E). G' – приведенная КС-грамматика.

3. Оформить работу согласно шагам.

Лабораторная работа №7-8

Построить МП-автомат Р и расширенный МП-автомат по КС-грамматики $G = (T, V, P, S)$, 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 Р и расширенного МП-автомата по заданной КС-грамматике (см. раздел 3.4.). Ответьте на вопрос, чем отличается МП-автомат Р 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());

```

Б). Определить последовательность тактов МП-автомата для выделенной цепочки.

2.2. А). Построить расширенный МП-автомат, используя алгоритм 3.9.

Б). 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 которых 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"
#pragma once
#include <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 M 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 M 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 таблица 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. $L = \{01 - (10)^* + \omega_1 01 \mid \omega_1 \in \{0,1\}^*\}$
3. $L = \{0(00)^* + 01\omega_1 \mid \omega_1 \in (0,1,2)^*\}$
4. $L = \{\omega_1\omega_2 1 \mid \omega_1 \in \{1,0\}^+, \omega_2 \in \{1,0\}^+\}$
5. $L = \{1\omega_1 1 - (00)^* \mid \omega_1 \in \{1,0\}^*\}$
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. $L = \{0(01)^* + 1\omega_1 0 \mid \omega_1 \in \{1,0\}^+\}$
16. $L = \{1\omega_1 1\omega_2 1 \mid \omega_1 \in \{0,1\}^+, \omega_2 \in \{0,1\}^+\}$
17. $L = \{011\omega_1 1(0)^* \mid \omega_1 \in \{0,1\}^+\}$
18. $L = \{\omega_1 0 - \omega_2 \mid \omega_1 \in \{0,1\}^+, \omega_2 \in \{0,1\}^+\}$
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 4-6

A). 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. $P = \{S \rightarrow b, S \rightarrow CF, S \rightarrow cCB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, C \rightarrow Ca, F \rightarrow d\}$

B). УстраниТЬ из грамматика $G = (T, V, P, S)$ 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. $P = \{S \rightarrow b, S \rightarrow C, S \rightarrow cCB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, C \rightarrow Ca, C \rightarrow \epsilon\}$

C). 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. $P = \{S \rightarrow b, S \rightarrow C, S \rightarrow cCB, A \rightarrow Ab, A \rightarrow c, B \rightarrow cB, C \rightarrow Ca, C \rightarrow b\}$

D). 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. $T = \{i, ^, -, (),\}, V = \{S, F, L\}, P = \{S \rightarrow (F ^ L), F \rightarrow - L, F \rightarrow i, L \rightarrow F\}$
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. $T = \{i, \&, ^, (),\}, V = \{S, F, L\}, P = \{S \rightarrow F^L, S \rightarrow (F), F \rightarrow \&L, F \rightarrow i, L \rightarrow F\}$
9. $T = \{i, \&, ^, (),\}, V = \{S, F, L\}, P = \{S \rightarrow F^L, S \rightarrow (S), F \rightarrow \&L, F \rightarrow i, L \rightarrow F\}$
10. $T = \{i, \&, ^, (),\}, V = \{S, F, L\}, P = \{S \rightarrow (F^L), S \rightarrow (S), F \rightarrow \&L, F \rightarrow i, L \rightarrow F\}$
11. $T = \{i, \&, ^, (),\}, V = \{S, F, L\}, P = \{S \rightarrow (F^L), F \rightarrow \&L, F \rightarrow i, L \rightarrow F\}$
12. $T = \{i, \&, ^, (),\}, V = \{S, F, L\}, P = \{S \rightarrow \&F^L, S \rightarrow (L), F \rightarrow L, F \rightarrow i, L \rightarrow F\}$
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. Axo A., 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 // РАН. Информационные технологии и Вычислительные системы. – М.: N2, 2005 стр.10-17.